

The Leray-Gårding method for finite difference schemes

Jean-François COULOMBEL*

May 6, 2015

Abstract

In [Ler53] and [Går56], LERAY and GÅRDING have developed a multiplier technique for deriving a priori estimates for solutions to scalar hyperbolic equations in either the whole space or the torus. In particular, the arguments in [Ler53, Går56] provide with at least one *local* multiplier and one *local* energy functional that is controlled along the evolution. The existence of such a local multiplier is the starting point of the argument by RAUCH in [Rau72] for the derivation of semigroup estimates for hyperbolic initial boundary value problems. In this article, we explain how this multiplier technique can be adapted to the framework of finite difference approximations of transport equations. The technique applies to numerical schemes with arbitrarily many time levels, and encompasses a somehow magical trick that has been known for a long time for the leap-frog scheme. More importantly, the existence and properties of the local multiplier enable us to derive *optimal* semigroup estimates for fully discrete hyperbolic initial boundary value problems, which answers a problem raised by TREFETHEN, KREISS and WU [Tre84, KW93].

AMS classification: 65M06, 65M12, 35L03, 35L04.

Keywords: hyperbolic equations, difference approximations, stability, boundary conditions, semigroup.

Throughout this article, we use the notation

$$\begin{aligned}\mathcal{U} &:= \{\zeta \in \mathbb{C}, |\zeta| > 1\}, \quad \overline{\mathcal{U}} := \{\zeta \in \mathbb{C}, |\zeta| \geq 1\}, \\ \mathbb{D} &:= \{\zeta \in \mathbb{C}, |\zeta| < 1\}, \quad \mathbb{S}^1 := \{\zeta \in \mathbb{C}, |\zeta| = 1\}, \quad \overline{\mathbb{D}} := \mathbb{D} \cup \mathbb{S}^1.\end{aligned}$$

We let $\mathcal{M}_n(\mathbb{K})$ denote the set of $n \times N$ matrices with entries in $\mathbb{K} = \mathbb{R}$ or \mathbb{C} . If $M \in \mathcal{M}_n(\mathbb{C})$, M^* denotes the conjugate transpose of M . We let I denote the identity matrix or the identity operator when it acts on an infinite dimensional space. We use the same notation x^*y for the Hermitian product of two vectors $x, y \in \mathbb{C}^n$ and for the Euclidean product of two vectors $x, y \in \mathbb{R}^n$. The norm of a vector $x \in \mathbb{C}^n$ is $|x| := (x^*x)^{1/2}$. The induced matrix norm on $\mathcal{M}_n(\mathbb{C})$ is denoted $\|\cdot\|$.

The letter C denotes a constant that may vary from line to line or within the same line. The dependence of the constants on the various parameters is made precise throughout the text.

*CNRS and Université de Nantes, Laboratoire de Mathématiques Jean Leray (UMR CNRS 6629), 2 rue de la Houssinière, BP 92208, 44322 Nantes Cedex 3, France. Email: jean-francois.coulombel@univ-nantes.fr. Research of the author was supported by ANR project BoND, ANR-13-BS01-0009-01.

In what follows, we let $d \geq 1$ denote a fixed integer, which will stand for the dimension of the space domain we are considering. We shall also use the space ℓ^2 of square integrable sequences. Sequences may be valued in \mathbb{C}^k for some integer k . Some sequences will be indexed by \mathbb{Z}^{d-1} while some will be indexed by \mathbb{Z}^d or a subset of \mathbb{Z}^d . We thus introduce some specific notation for the norms. Let $\Delta x_i > 0$ for $i = 1, \dots, d$ be d space steps. We shall make use of the $\ell^2(\mathbb{Z}^{d-1})$ -norm that we define as follows: for all $v \in \ell^2(\mathbb{Z}^{d-1})$,

$$\|v\|_{\ell^2(\mathbb{Z}^{d-1})}^2 := \left(\prod_{k=2}^d \Delta x_k \right) \sum_{i=2}^d \sum_{j_i \in \mathbb{Z}} |v_{j_2, \dots, j_d}|^2.$$

The corresponding scalar product is denoted $\langle \cdot, \cdot \rangle_{\ell^2(\mathbb{Z}^{d-1})}$. Then for all integers $m_1 \leq m_2$, we set

$$\|u\|_{m_1, m_2}^2 := \Delta x_1 \sum_{j_1=m_1}^{m_2} \|u_{j_1, \cdot}\|_{\ell^2(\mathbb{Z}^{d-1})}^2,$$

to denote the ℓ^2 -norm on the set $[m_1, m_2] \times \mathbb{Z}^{d-1}$ (m_1 may equal $-\infty$ and m_2 may equal $+\infty$). The corresponding scalar product is denoted $\langle \cdot, \cdot \rangle_{m_1, m_2}$. Other notation throughout the text is meant to be self-explanatory.

1 Introduction

1.1 Some motivations and a brief reminder

The ultimate goal of this article is to derive semigroup estimates for finite difference approximations of hyperbolic initial boundary value problems. Up to now, the only available general stability theory for such numerical schemes is due to GUSTAFSSON, KREISS and SUNDSTRÖM [GKS72]. It relies on a Laplace transform with respect to the time variable, and the corresponding stability estimates are thereby restricted to zero initial data. A long standing problem in this line of research is, starting from the GKS stability estimates, which are *resolvent* type estimates, to incorporate nonzero initial data and to derive *semigroup* estimates, see, e.g., the discussion in [Tre84, section 4]. This problem is delicate for the following reason: the validity of the GKS stability estimate is known to be equivalent to a *slightly stronger version* of the resolvent estimate

$$\sup_{z \in \mathcal{U}} (|z| - 1) \|(zI - T)^{-1}\|_{\mathcal{L}(\ell^2(\mathbb{N}))} < +\infty, \quad (1)$$

where T is some bounded operator on $\ell^2(\mathbb{N})$ that incorporates both the discretization of the hyperbolic equation and the numerical boundary conditions. Deriving an optimal semigroup estimate amounts to showing that T is power bounded. In finite dimension, the equivalence between power boundedness of T and the resolvent condition (1) is known as the KREISS matrix Theorem, but the analogous equivalence is known to fail in general in infinite dimension. Worse, even the strong resolvent condition

$$\sup_{n \geq 1} \sup_{z \in \mathcal{U}} (|z| - 1)^n \|(zI - T)^{-n}\|_{\mathcal{L}(\ell^2(\mathbb{N}))} < +\infty,$$

does not imply in general that T is power bounded, see, e.g., the review [SW97] or [TE05] for details and historical comments.

Optimal semigroup estimates have nevertheless been derived for some discretized hyperbolic initial boundary value problems. More specifically, the first general derivation of semigroup estimates is due

to WU [Wu95], whose analysis deals with numerical schemes with two time levels and scalar equations. The results in [Wu95] were extended by GLORIA and the author in [CG11] to systems in arbitrary space dimension, but the arguments in [CG11] are still restricted to numerical schemes with two time levels. The present article gives, as far as we are aware of, the first systematic derivation of semigroup estimates for fully discrete hyperbolic initial boundary value problems in the case of numerical schemes with arbitrarily many time levels. It generalizes the arguments of [Wu95, CG11] and provides new insight for the construction of “dissipative” numerical boundary conditions for discretized evolution equations. Let us observe that the leap-frog scheme, with some specific boundary conditions, has been dealt with by THOMAS [Tho72] by using a *multiplier* technique. It is precisely this technique which we aim at developing in a systematic fashion for numerical schemes with arbitrarily many time levels. In particular, we shall explain why the somehow magical multiplier $u_j^{n+2} + u_j^n$ for the leap-frog scheme, see, e.g., [RM67], follows from a general theory that is the analogue of the LERAY-GÅRDING method for partial differential equations, which we briefly recall now.

The method by LERAY and GÅRDING [Ler53, Går56] provides with suitable multipliers for scalar hyperbolic operators of arbitrary order. Namely, given an integer $m \geq 0$, we consider a partial differential operator of the form

$$L := \partial_t^{m+1} + \sum_{k=1}^{m+1} P_k(\partial_x) \partial_t^{m+1-k},$$

where $t \in \mathbb{R}$ stands for the time variable, $x \in \mathbb{R}^d$ stands for the space variable¹, and each operator $P_k(\partial_x)$ is a linear combination of spatial partial derivatives of order k :

$$P_k(\partial_x) = \sum_{|\alpha|=k} p_{k,\alpha} \partial_x^\alpha, \quad \partial_x^\alpha := \partial_{x_1}^{\alpha_1} \cdots \partial_{x_d}^{\alpha_d}, \quad |\alpha| := \alpha_1 + \cdots + \alpha_d.$$

In the above formula, the $p_{k,\alpha}$ ’s are real numbers². Well-posedness of the Cauchy problem

$$L u = 0, \quad (u, \partial_t u, \dots, \partial_t^m u)|_{t=0} = (u_0, u_1, \dots, u_m), \quad (2)$$

in Sobolev spaces is known to be linked with *hyperbolicity* of L . Namely, if L is strictly hyperbolic, meaning that for all $\xi \in \mathbb{R}^d \setminus \{0\}$, the (homogeneous) polynomial

$$P(\tau, \xi) := \tau^{m+1} + \sum_{k=1}^{m+1} P_k(i\xi) \tau^{m+1-k}, \quad P_k(i\xi) := i^k \sum_{|\alpha|=k} p_{k,\alpha} \xi^\alpha, \quad (3)$$

has $m+1$ simple purely imaginary roots with respect to τ , then the Cauchy problem (2) is well-posed in $H^m(\mathbb{R}^d) \times \cdots \times L^2(\mathbb{R}^d)$. In particular, there exists a constant $C > 0$, that is independent of the solution u and the initial data u_0, u_1, \dots, u_m , such that there holds:

$$\sup_{t \in \mathbb{R}} \sum_{k=0}^m \|\partial_t^k u(t)\|_{H^{m-k}(\mathbb{R}^d)} \leq C \sum_{k=0}^m \|u_k\|_{H^{m-k}(\mathbb{R}^d)}. \quad (4)$$

The method by LERAY and GÅRDING gives a quick and elegant way to derive the estimate (4) assuming that the solution u to (2) is sufficiently smooth. By standard duality arguments, the validity of the a

¹The periodic case $x \in \mathbb{T}^d$ can be dealt with in a similar way and is actually the one considered in [Går56].

²We restrict here for simplicity to linear operators with constant coefficients.

priori estimate (4) yields well-posedness -meaning existence, uniqueness and continuous dependence on the data- for (2). Hence the main point is to prove (4) assuming that u is sufficiently smooth and decaying at infinity so that all integration by parts arising in the computations are legitimate. The main idea is to find a suitable quantity Mu , which we call a multiplier and that will be linear with respect to u , such that when integrating the quantity $0 = (Mu)(Lu)$ on the slab $[0, T] \times \mathbb{R}^d$, one gets the estimate (4) *for free* (negative times are obtained by changing $t \rightarrow -t$). Following [Ler53, Chapter VI] and [Går56, Section 3], one possible choice of a multiplier is given by $L'u$ where L' stands for the partial differential operator of order m whose symbol is $\partial_\tau P$, with P given in (3). Why $L'u$ is a good multiplier is justified in [Ler53, Går56]. A well-known particular case is the choice of $2\partial_t u$ as a multiplier for the wave equation. Here $P(\tau, \xi) = \tau^2 + |\xi|^2$ and therefore $\partial_\tau P = 2\tau$, hence the choice $2\partial_t u$. The latter quantity is indeed a suitable multiplier for the wave operator because of the formula³:

$$2\partial_t u (\partial_t^2 u - \Delta_x u) = \partial_t \left((\partial_t u)^2 + \sum_{j=1}^d (\partial_{x_j} u)^2 \right) - 2 \operatorname{div}_x (\partial_t u \nabla_x u).$$

The important fact here is that the *energy*:

$$(\partial_t u)^2 + \sum_{j=1}^d (\partial_{x_j} u)^2,$$

is a positive definite quadratic form of the first order partial derivatives of u . Let us observe that the multiplier $L'u$ is *local*, meaning that its pointwise value at (t, x) only depends on u in a neighborhood of (t, x) . This is important in view of using this multiplier in the study of initial boundary value problems. Another important remark is that the above energy is also *local*, and the arguments in [Ler53, Går56] show that this property is not specific to the wave operator. The fact that both the multiplier and the energy are local is crucial in the arguments of [Rau72, Lemma 1]. In our framework of discretized equations, the multiplier will be local but the energy will not necessarily be so. We shall not exactly follow the arguments of [Rau72] which use time reversibility, but rather construct dissipative boundary conditions which will yield the optimal semigroup estimate we are aiming at.

1.2 The main result

We first set a few notations. We let $\Delta x_1, \dots, \Delta x_d, \Delta t > 0$ denote space and time steps where the ratios, the so-called COURANT-FRIEDRICHS-LEWY parameters, $\lambda_i := \Delta t / \Delta x_i$, $i = 1, \dots, d$, are fixed positive constants. We keep $\Delta t \in (0, 1]$ as a small parameter and let the space steps $\Delta x_1, \dots, \Delta x_d$ vary accordingly. The ℓ^2 -norms with respect to the space variables have been previously defined and thus depend on Δt and the CFL parameters through the mesh volume ($\Delta x_2 \cdots \Delta x_d$ on \mathbb{Z}^{d-1} , and $\Delta x_1 \cdots \Delta x_d$ on \mathbb{Z}^d). We always identify a sequence w indexed by either \mathbb{N} (for time), \mathbb{Z}^{d-1} or \mathbb{Z}^d (for space), with the corresponding step function. In particular, we shall feel free to take Fourier or Laplace transforms of such sequences.

For all $j \in \mathbb{Z}^d$, we set $j = (j_1, j')$ with $j' := (j_2, \dots, j_d) \in \mathbb{Z}^{d-1}$. We let $p, q, r \in \mathbb{N}^d$ denote some fixed multi-integers, and define $p_1, q_1, r_1, p', q', r'$ according to the above notation. We also let $s \in \mathbb{N}$ denote

³We refer to [Går56, page 74] for the generalization of such "integration by parts" formula to partial derivatives of higher order.

some fixed integer. We consider a recurrence relation of the form:

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_{\sigma} u_j^{n+\sigma} = \Delta t F_j^{n+s+1}, & j \in \mathbb{Z}^d, \quad j_1 \geq 1, \quad n \geq 0, \\ u_j^{n+s+1} + \sum_{\sigma=0}^{s+1} B_{j_1, \sigma} u_{1, j'}^{n+\sigma} = g_j^{n+s+1}, & j \in \mathbb{Z}^d, \quad j_1 = 1 - r_1, \dots, 0, \quad n \geq 0, \\ u_j^n = f_j^n, & j \in \mathbb{Z}^d, \quad j_1 \geq 1 - r_1, \quad n = 0, \dots, s, \end{cases} \quad (5)$$

where the operators Q_{σ} and $B_{j_1, \sigma}$ are given by:

$$Q_{\sigma} := \sum_{\ell_1=-r_1}^{p_1} \sum_{\ell'=-r'}^{p'} a_{\ell, \sigma} \mathbf{S}^{\ell}, \quad B_{j_1, \sigma} := \sum_{\ell_1=0}^{q_1} \sum_{\ell'=-q'}^{q'} b_{\ell, j_1, \sigma} \mathbf{S}^{\ell}. \quad (6)$$

In (6), the $a_{\ell, \sigma}, b_{\ell, j_1, \sigma}$ are *real numbers* and are independent of the small parameter Δt (they may depend on the CFL parameters though), while \mathbf{S} denotes the shift operator on the space grid: $(\mathbf{S}^{\ell} v)_j := v_{j+\ell}$ for $j, \ell \in \mathbb{Z}^d$. We have also used the short notation

$$\sum_{\ell'=-r'}^{p'} := \sum_{i=2}^d \sum_{\ell_i=-r_i}^{p_i}, \quad \sum_{\ell'=-q'}^{q'} := \sum_{i=2}^d \sum_{\ell_i=-q_i}^{q_i}.$$

The numerical scheme (5) is understood as follows: one starts with ℓ^2 initial data $(f_j^0), \dots, (f_j^s)$ defined for $j_1 \geq 1 - r_1$. Assuming that the solution has been defined up to some time index $n + s$, $n \geq 0$, then the first and second equations in (5) should uniquely determine u_j^{n+s+1} for $j_1 \geq 1 - r_1$, $j' \in \mathbb{Z}^{d-1}$. The meshes associated with $j_1 \geq 1$ correspond to the *interior domain* while those associated with $j_1 = 1 - r_1, \dots, 0$ represent the *discrete boundary*. We wish to deal here simultaneously with explicit and implicit schemes and therefore make the following solvability assumption.

Assumption 1 (Solvability of (5)). *The operator Q_{s+1} is an isomorphism on $\ell^2(\mathbb{Z}^d)$. Moreover, for all $(F_j) \in \ell^2(\mathbb{N}^* \times \mathbb{Z}^{d-1})$ and for all $g_{1-r_1}, \dots, g_0 \in \ell^2(\mathbb{Z}^{d-1})$, there exists a unique solution $(u_j)_{j_1 \geq 1-r_1} \in \ell^2$ to the equations*

$$\begin{cases} Q_{s+1} u_j = F_j, & j \in \mathbb{Z}^d, \quad j_1 \geq 1, \\ u_j + B_{j_1, s+1} u_{1, j'} = g_j, & j \in \mathbb{Z}^d, \quad j_1 = 1 - r_1, \dots, 0. \end{cases}$$

In particular, Assumption 1 is trivially satisfied in the case of explicit schemes for which Q_{s+1} is the identity ($a_{\ell, s+1} = \delta_{\ell_1, 0} \cdots \delta_{\ell_d, 0}$ in (6), with δ the Kronecker symbol).

The first and second equations in (5) therefore uniquely determine u_j^{n+s+1} for $j_1 \geq 1 - r_1$, and one then proceeds to the following time index $n + s + 2$. Existence and uniqueness of a solution (u_j^n) to (5) follows from Assumption 1, so the last requirement for well-posedness is continuous dependence of the solution on the three possible source terms $(F_j^n), (g_j^n), (f_j^n)$. This is a *stability* problem for which several definitions can be chosen according to the functional framework. The following one dates back to [GKS72] in one space dimension and was also considered by MICHELSON [Mic83] in several space dimensions. It is specifically relevant when the boundary conditions are non-homogeneous ($(g_j^n) \neq 0$):

Definition 1 (Strong stability). *The finite difference approximation (5) is said to be "strongly stable" if there exists a constant C such that for all $\gamma > 0$ and all $\Delta t \in (0, 1]$, the solution (u_j^n) to (5) with*

$(f_j^0) = \dots = (f_j^s) = 0$ satisfies the estimate:

$$\begin{aligned} & \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j=1-r_1}^{p_1} \|u_{j1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ & \leq C \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (7) \end{aligned}$$

The main contributions in [GKS72, Mic83] are to show that strong stability can be characterized by a certain *algebraic condition*, which is usually referred to as the Uniform KREISS-LOPATINSKII Condition, see [Cou13] for an overview of such results. We do not pursue such arguments here but rather assume from the start that (5) is strongly stable. We can thus control, with zero initial data, ℓ^2 type norms of the solution to (5). Our goal is to understand which kind of stability estimate holds for the solution to (5) when one now considers nonzero initial data $(f_j^0), \dots, (f_j^s)$ in ℓ^2 . Our main assumption is the following.

Assumption 2 (Stability for the discrete Cauchy problem). *For all $\xi \in \mathbb{R}^d$, the dispersion relation*

$$\sum_{\sigma=0}^{s+1} \widehat{Q}_\sigma(e^{i\xi_1}, \dots, e^{i\xi_d}) z^\sigma = 0, \quad \widehat{Q}_\sigma(\kappa) := \sum_{\ell=-r}^p \kappa^\ell a_{\ell, \sigma}, \quad (8)$$

has $s+1$ simple roots in $\overline{\mathbb{D}}$. (The von Neumann condition is said to hold when the roots are located in $\overline{\mathbb{D}}$.) In (8), we have used the classical notation

$$\kappa^\ell := \kappa_1^{\ell_1} \dots \kappa_d^{\ell_d},$$

for $\kappa \in (\mathbb{C} \setminus \{0\})^d$ and $\ell \in \mathbb{Z}^d$.

From Assumption 1, we know that Q_{s+1} is an isomorphism on ℓ^2 , which implies by Fourier analysis that $\widehat{Q}_{s+1}(e^{i\xi_1}, \dots, e^{i\xi_d})$ does not vanish for any $\xi \in \mathbb{R}^d$. In particular, the dispersion relation (8) is a polynomial equation of degree $s+1$ in z for any $\xi \in \mathbb{R}^d$. We now make the following assumption, which already appeared in [GKS72, Mic83] and several other works on the same topic.

Assumption 3 (Noncharacteristic discrete boundary). *For $\ell_1 = -r_1, \dots, p_1$, $z \in \mathbb{C}$ and $\eta \in \mathbb{R}^{d-1}$, let us define*

$$a_{\ell_1}(z, \eta) := \sum_{\sigma=0}^{s+1} z^\sigma \sum_{\ell'=-r'}^{p'} a_{\ell, \sigma} e^{i\ell' \cdot \eta}. \quad (9)$$

Then a_{-r_1} and a_{p_1} do not vanish on $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$, and they have nonzero degree with respect to z for all $\eta \in \mathbb{R}^{d-1}$.

Our main result is comparable with [Wu95, Theorem 3.3] and [CG11, Theorems 2.4 and 3.5] and shows that strong stability (or "GKS stability") is a sufficient condition for incorporating ℓ^2 initial conditions in (5) and proving *optimal* semigroup estimates. The main price to pay in Assumption 2 is that the roots of the dispersion relation (8), which are nothing but the eigenvalues of the so-called *amplification matrix* for the Cauchy problem, need to be *simple*. This property is satisfied for instance by the leap-frog and modified leap-frog schemes in several space dimensions, under an appropriate CFL condition, see Paragraph 1.3. Our main result reads as follows.

Theorem 1. *Let Assumptions 1, 2 and 3 be satisfied, and assume that the scheme (5) is strongly stable in the sense of Definition 1. Then there exists a constant C such that for all $\gamma > 0$ and all $\Delta t \in (0, 1]$, the solution to (5) satisfies the estimate:*

$$\begin{aligned} & \sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 + \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 \\ & + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2 + \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 \right. \\ & \left. + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (10) \end{aligned}$$

In particular, the scheme (5) is "semigroup stable" in the sense that there exists a constant C such that for all $\Delta t \in (0, 1]$, the solution (u_j^n) to (5) with $(F_j^n) = (g_j^n) = 0$ satisfies the estimate

$$\sup_{n \geq 0} \|u^n\|_{1-r_1, +\infty}^2 \leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2. \quad (11)$$

The scheme (5) is also ℓ^2 -stable with respect to boundary data, see [Tre84, Definition 4.5], in the sense that there exists a constant C such that for all $\Delta t \in (0, 1]$, the solution (u_j^n) to (5) with $(F_j^n) = (f_j^n) = 0$ satisfies the estimate

$$\sup_{n \geq 0} \|u^n\|_{1-r_1, +\infty}^2 \leq C \sum_{n \geq s+1} \Delta t \sum_{j_1=1-r_1}^0 \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2.$$

Theorem 1 gives the optimal semigroup estimate (11), and is therefore an improvement with respect to our earlier work [Cou14] where in one space dimension, and under an appropriate *non-glancing* condition⁴, we were able to derive the estimate (here $r_1 = r$, $p_1 = p$ since $d = 1$):

$$\begin{aligned} & \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r, +\infty}^2 + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j=1-r}^p |u_j^n|^2 \\ & \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r, +\infty}^2 + \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j=1-r}^0 |g_j^n|^2 \right\}. \end{aligned}$$

The latter estimate does not incorporate on the left hand side the quantity:

$$\sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{1-r, +\infty}^2,$$

and was unfortunately still not sufficient for deriving the semigroup estimate (11). Our main contribution in this article is to exhibit a suitable multiplier for the multistep recurrence relation in (5). With this multiplier, we can readily show that, for zero initial data, the (discrete) derivative of an *energy* can be controlled, as in [Rau72], by the trace estimate of (u_j^n) and this is where strong stability comes into play.

⁴The non-glancing condition is unfortunately not met by the leap-frog scheme.

This first argument gives Theorem 1 for zero initial data (and even for nonzero initial data if the non-glancing condition of [Cou14] is satisfied). By linearity we can then reduce to the case of zero forcing terms in the interior and on the boundary. The next arguments in [Rau72] use time reversibility, which basically always fails for numerical schemes⁵. Hence we must find another argument for dealing with nonzero initial data. Hopefully, the properties of our multiplier enable us to construct an auxiliary problem, where we modify the boundary conditions of (5), and for which we can prove optimal semigroup and trace estimates by "hand-made" calculations. In other words, we exhibit an alternative set of boundary conditions that yields *strict dissipativity*. Using these auxiliary numerical boundary conditions, the proof of Theorem 1 follows from a standard superposition argument, see, e.g., [BGS07, Section 4.5] for partial differential equations or [Wu95, CG11] for numerical schemes.

Remark 1. *Assumption 3 excludes the case of explicit two level schemes for which $s = 0$ and $Q_1 = I$, for in that case a_{-r_1} and/or a_{p_1} do not depend on z . However, this case has already been dealt with in [Wu95, CG11], and we shall see in Section 3 where the assumption that a_{-r_1} and a_{p_1} are not constant is involved, and why the proof is actually simpler in the case $s = 0$ and $Q_1 = I$.*

1.3 Examples

1.3.1 One space dimension

Our goal is to approximate the outgoing transport equation ($d = 1$ here):

$$\partial_t u + a \partial_x u = 0, \quad u|_{t=0} = u_0, \quad (12)$$

with $t, x > 0$ and $a < 0$. The latter transport equation does not require any boundary condition at $x = 0$. However, discretizing (12) usually requires prescribing numerical boundary conditions, unless one considers an upwind type scheme with a space stencil "on the right" (meaning $r_1 = 0$ in (5)). We now detail two possible multistep schemes for discretizing (12). Both are obtained by the so-called method of lines, which amounts to first discretizing the space derivative $\partial_x u$ and then choosing an integration technique for discretizing the time evolution, see [GKO95].

The leap-frog scheme. It is obtained by approximating the space derivative $\partial_x u$ by the centered difference $(u_{j+1} - u_{j-1})/(2\Delta x)$, and by then applying the so-called Nyström method of order 2, see [HNW93, Chapter III.1]. The resulting approximation reads

$$u_j^{n+2} + \lambda a (u_{j+1}^{n+1} - u_{j-1}^{n+1}) - u_j^n = 0,$$

which corresponds to $s = p = r = 1$. Recall that $\lambda > 0$ denotes the fixed ratio $\Delta t/\Delta x$. Even though (12) does not require any boundary condition at $x = 0$, the leap-frog scheme stencil includes one point to the left, and we therefore need to prescribe some numerical boundary condition at $j = 0$. One possibility⁶ is to prescribe the homogeneous or inhomogeneous Dirichlet boundary condition. With general source terms, the corresponding scheme reads

$$\begin{cases} u_j^{n+2} + \lambda a (u_{j+1}^{n+1} - u_{j-1}^{n+1}) - u_j^n = \Delta t F_j^{n+2}, & j \geq 1, \quad n \geq 0, \\ u_0^{n+2} = g_0^{n+2}, & n \geq 0, \\ (u_j^0, u_j^1) = (f_j^0, f_j^1), & j \geq 0. \end{cases} \quad (13)$$

⁵With the notable exception of the leap-frog scheme that is indeed time reversible !

⁶This is of course not the only possibility and we refer to [GKO95, Oli74, Slo83, Tre84] for some other possible choices which might be more meaningful from a consistency and accuracy point of view. Our main concern here is a discussion on stability for (5) and the Dirichlet boundary conditions are a good illustration for this aspect.

Assumption 1 is trivially satisfied because (13) is explicit. The leap-frog scheme satisfies Assumption 2 provided that $\lambda|a| < 1$. In that case, the two roots to the dispersion relation

$$z^2 + 2i\lambda a \sin \xi z - 1 = 0,$$

are simple and have modulus 1 for all $\xi \in \mathbb{R}$. Assumption 3 is satisfied as long as the velocity a is nonzero, for in that case $a_1(z) = -a_{-1}(z) = \lambda a z$. The scheme (13) is known to be strongly stable, see [GT81]. In particular, Theorem 1 shows that (13) is semigroup stable. An illustration of this stability property is given in the numerical simulation of a bump function, propagating at speed $a = -1$ towards the left. Homogeneous Dirichlet boundary conditions are enforced at $j = 0$. The reflection of the bump generates a highly oscillatory wave packet that propagates with velocity $+1$ towards the right. The envelope of this wave packet coincides with the profile of the initial condition, which indicates that the ℓ^2 -norm is roughly preserved by the evolution. This numerical observation is in agreement with semigroup boundedness.

Other choices of numerical boundary conditions for the leap-frog scheme or its fourth order extension are discussed, e.g., in [Oli74, Slo83, Tho72, Tre84]. The main discussion in [Oli74, Slo83, Tre84] is to verify strong stability for a wide choice of numerical boundary conditions, and if strong stability holds, then Theorem 1 automatically gives semigroup boundedness, which was not achieved in these earlier works.

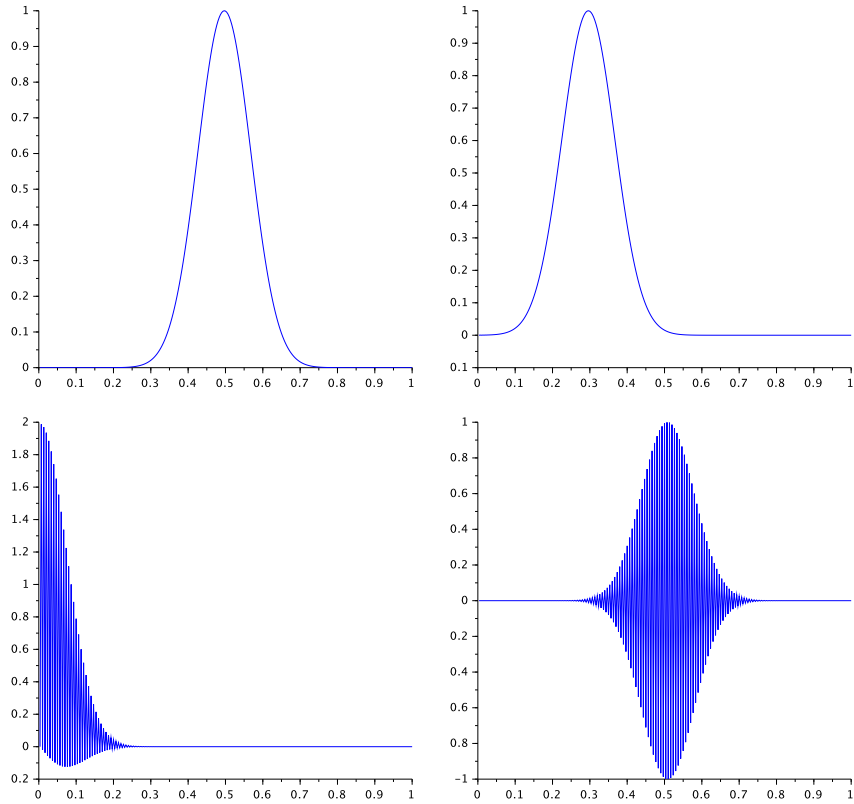


Figure 1: Reflection of a bump by the leap-frog scheme with homogeneous Dirichlet condition at four successive times.

A scheme based on the backwards differentiation rule. We still start from the transport equation (12), approximate the space derivative $\partial_x u$ by the centered finite difference $(u_{j+1} - u_{j-1})/(2 \Delta x)$, and then apply the backwards differentiation formula of order 2, see [HNW93, Chapter III.1]. The resulting scheme reads:

$$\frac{3}{2} u_j^{n+2} + \frac{\lambda a}{2} (u_{j+1}^{n+2} - u_{j-1}^{n+2}) - 2 u_j^{n+1} + \frac{1}{2} u_j^n = 0.$$

This corresponds to $s = 1$ and

$$Q_2 : (u_j)_{j \in \mathbb{Z}} \mapsto \left(\frac{3}{2} u_j + \frac{\lambda a}{2} (u_{j+1} - u_{j-1}) \right)_{j \in \mathbb{Z}}.$$

The operator Q_2 is an isomorphism on $\ell^2(\mathbb{Z})$ since Q_2 is an isomorphism for any small λa (as a perturbation of $3/2 I$), Q_2 depends continuously on λa , and there holds (uniformly with respect to λa):

$$\frac{3}{2} \|u\|_{-\infty, +\infty} \leq \|Q_2 u\|_{-\infty, +\infty}.$$

The operator Q_2 is therefore an isomorphism on $\ell^2(\mathbb{Z})$ for any $\lambda a > 0$ (see, e.g., [Cou09, Lemma 4.3]). Let us now study the dispersion relation (8), which reads here

$$\left(\frac{3}{2} + i \lambda a \sin \xi \right) z^2 - 2z + \frac{1}{2} = 0.$$

It is clear that the latter equation has two simple roots in z for any $\xi \in \mathbb{R}$. Moreover, if $\sin \xi = 0$, the roots are 1 and 1/3 which belong to $\overline{\mathbb{D}}$. In the case $\sin \xi \neq 0$, none of the roots belongs to \mathbb{S}^1 and examining the case $\lambda a \sin \xi = 1$, we find that for $\sin \xi \neq 0$, both roots belong to \mathbb{D} (which is consistent with the shape of the stability region for the backwards differentiation formula of order 2, see [HW96, Chapter V.1]). Assumption 2 is therefore satisfied. Assumption 3 is satisfied as long as a is nonzero since there holds $p = r = 1$ and $a_1(z) = a_{-1}(z) = \lambda a z^2/2$.

Theorem 1 therefore yields semigroup boundedness as long as one uses numerical boundary conditions for which the numerical scheme is well-defined (this is at least the case for λa small enough) and strong stability holds.

1.3.2 Two space dimensions

Here we wish to approximate the two-dimensional transport equation ($d = 2$):

$$\partial_t u + a_1 \partial_{x_1} u + a_2 \partial_{x_2} u = 0, \quad u|_{t=0} = u_0,$$

in the space domain $\{x_1 > 0, x_2 \in \mathbb{R}\}$. When a_1 is negative, the latter problem does not necessitate any boundary condition at $x_1 = 0$. Following [AG76], we use one of the following two-dimensional versions of the leap-frog scheme, either

$$u_{j,k}^{n+2} + \lambda_1 a_1 (u_{j+1,k}^{n+1} - u_{j-1,k}^{n+1}) + \lambda_2 a_2 (u_{j,k+1}^{n+1} - u_{j,k-1}^{n+1}) - u_{j,k}^n = 0, \quad (14)$$

or

$$\begin{aligned} u_{j,k}^{n+2} + \lambda_1 a_1 \left(\frac{u_{j+1,k+1}^{n+1} + u_{j+1,k-1}^{n+1}}{2} - \frac{u_{j-1,k+1}^{n+1} + u_{j-1,k-1}^{n+1}}{2} \right) \\ + \lambda_2 a_2 \left(\frac{u_{j+1,k+1}^{n+1} + u_{j-1,k+1}^{n+1}}{2} - \frac{u_{j+1,k-1}^{n+1} + u_{j-1,k-1}^{n+1}}{2} \right) - u_{j,k}^n = 0. \end{aligned} \quad (15)$$

Assumption 1 is trivially satisfied because (14) and (15) are explicit schemes. The scheme (14) satisfies Assumption 2 if and only if $\lambda_1 |a_1| + \lambda_2 |a_2| < 1$, while the scheme (15) satisfies Assumption 2 if and only if $\max(\lambda_1 |a_1|, \lambda_2 |a_2|) < 1$. Let us now study when Assumption 3 is valid. For the scheme (14), we have $r_1 = p_1 = 1$, and

$$a_1(z, \eta) = \lambda_1 a_1 z, \quad a_{-1}(z, \eta) = -a_1(z, \eta),$$

so Assumption 3 is valid as long as $a_1 \neq 0$. For the scheme (15), we have again $r_1 = p_1 = 1$, and

$$a_1(z, \eta) = z(\lambda_1 a_1 \cos \eta + i \lambda_2 a_2 \sin \eta), \quad a_{-1}(z, \eta) = z(-\lambda_1 a_1 \cos \eta + i \lambda_2 a_2 \sin \eta),$$

so Assumption 3 is valid as long as both a_1 and a_2 are nonzero. We refer to [AG79] for the verification of strong stability depending on the choice of some numerical boundary conditions for (14) or (15). Once again, if strong stability holds, then Theorem 1 yields semigroup boundedness and ℓ^2 -stability with respect to boundary data.

2 The Leray-Gårding method for fully discrete Cauchy problems

This section is devoted to proving stability estimates for discretized Cauchy problems, which is the first step before considering the discretized initial boundary value problem (5). More precisely, we consider the simpler case of the whole space $j \in \mathbb{Z}^d$, and the recurrence relation:

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma u_j^{n+\sigma} = 0, & j \in \mathbb{Z}^d, \quad n \geq 0, \\ u_j^n = f_j^n, & j \in \mathbb{Z}^d, \quad n = 0, \dots, s, \end{cases} \quad (16)$$

where the operators Q_σ are given by (6). We recall that in (6), the $a_{\ell, \sigma}$ are real numbers and are independent of the small parameter Δt (they may depend on the CFL parameters $\lambda_1, \dots, \lambda_d$), while \mathbf{S} denotes the shift operator on the space grid: $(\mathbf{S}^\ell v)_j := v_{j+\ell}$ for $j, \ell \in \mathbb{Z}^d$. Stability of (16) is defined as follows.

Definition 2 (Stability for the discrete Cauchy problem). *The numerical scheme defined by (16) is (ℓ^2) -stable if Q_{s+1} is an isomorphism from $\ell^2(\mathbb{Z}^d)$ onto itself, and if furthermore there exists a constant $C_0 > 0$ such that for all $\Delta t \in (0, 1]$, for all initial conditions $(f_j^0)_{j \in \mathbb{Z}^d}, \dots, (f_j^s)_{j \in \mathbb{Z}^d}$ in $\ell^2(\mathbb{Z}^d)$, there holds*

$$\sup_{n \in \mathbb{N}} \|u^n\|_{-\infty, +\infty}^2 \leq C_0 \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2. \quad (17)$$

Let us quickly recall that stability in the sense of Definition 2 is in fact independent of $\Delta t \in (0, 1]$ (because (16) does not involve Δt and (17) can be simplified on either side by $\prod_i \Delta x_i$), and can be characterized in terms of the uniform power boundedness of the so-called amplification matrix

$$\mathcal{A}(\kappa) := \begin{pmatrix} -\widehat{Q_s}(\kappa)/\widehat{Q_{s+1}}(\kappa) & \dots & \dots & -\widehat{Q_0}(\kappa)/\widehat{Q_{s+1}}(\kappa) \\ 1 & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ 0 & 0 & 1 & 0 \end{pmatrix} \in \mathcal{M}_{s+1}(\mathbb{C}), \quad (18)$$

where the $\widehat{Q_\sigma}(\kappa)$'s are defined in (8) and where it is understood that \mathcal{A} is defined on the largest open set of \mathbb{C}^d on which $\widehat{Q_{s+1}}$ does not vanish. Let us also recall that if Q_{s+1} is an isomorphism from $\ell^2(\mathbb{Z}^d)$

onto itself, then $\widehat{Q_{s+1}}$ does not vanish on $(\mathbb{S}^1)^d$, and therefore does not vanish on an open neighborhood of $(\mathbb{S}^1)^d$. With the above definition (18) for \mathcal{A} , the following well-known result holds:

Proposition 1 (Characterization of stability for the fully discrete Cauchy problem). *Assume that Q_{s+1} is an isomorphism from $\ell^2(\mathbb{Z}^d)$ onto itself. Then the scheme (16) is stable in the sense of Definition 2 if and only if there exists a constant $C_1 > 0$ such that the amplification matrix \mathcal{A} in (18) satisfies*

$$\forall n \in \mathbb{N}, \quad \forall \xi \in \mathbb{R}^d, \quad \left| \mathcal{A}(e^{i\xi_1}, \dots, e^{i\xi_d})^n \right| \leq C_1. \quad (19)$$

In particular, the spectral radius of $\mathcal{A}(e^{i\xi_1}, \dots, e^{i\xi_d})$ should not be larger than 1 (the so-called von Neumann condition).

The eigenvalues of $\mathcal{A}(e^{i\xi_1}, \dots, e^{i\xi_d})$ are the roots to the dispersion relation (8). When these roots are simple for all $\xi \in \mathbb{R}^d$, the von Neumann condition is both necessary and *sufficient* for stability of (16), see, e. g., [Cou13, Proposition 3]. Assumption 2 is therefore a way to assume that (16) is stable for the discrete Cauchy problem. Our goal is to derive the semigroup estimate (17) not by applying Fourier transform to (16) and using uniform power boundedness of \mathcal{A} , but rather by multiplying the first equation in (16) by a suitable *local* multiplier. The analysis relies first on the simpler case where one only considers the time evolution and no additional space variable.

2.1 Stable recurrence relations

In this Paragraph, we consider sequences $(v^n)_{n \in \mathbb{N}}$ with values in \mathbb{C} . The index n should be thought of as the discrete time variable, and we therefore introduce the new notation \mathbf{T} for the shift operator on the time grid: $(\mathbf{T}^m v)^n := v^{n+m}$ for all $m, n \in \mathbb{N}$. We start with the following elementary but crucial Lemma, which is the analogue of [Gär56, Lemme 1.1].

Lemma 1 (The energy-dissipation balance law). *Let $P \in \mathbb{C}[X]$ be a polynomial of degree $s+1$ whose roots are simple and located in $\overline{\mathbb{D}}$. Then there exists a positive definite Hermitian form q_e on \mathbb{C}^{s+1} , and a nonnegative Hermitian form q_d on \mathbb{C}^{s+1} , that both depend in a \mathcal{C}^∞ way on P , such that for any sequence $(v^n)_{n \in \mathbb{N}}$ with values in \mathbb{C} , there holds*

$$\forall n \in \mathbb{N}, \quad 2 \operatorname{Re} \left(\overline{\mathbf{T}(P'(\mathbf{T})v^n)} P(\mathbf{T})v^n \right) = (s+1) |P(\mathbf{T})v^n|^2 + (\mathbf{T} - I)(q_e(v^n, \dots, v^{n+s})) + q_d(v^n, \dots, v^{n+s}).$$

In particular, for all sequence $(v^n)_{n \in \mathbb{N}}$ that satisfies the recurrence relation

$$\forall n \in \mathbb{N}, \quad P(\mathbf{T})v^n = 0,$$

the sequence $(q_e(v^n, \dots, v^{n+s}))_{n \in \mathbb{N}}$ is nonincreasing.

The fact that there exists a Hermitian norm on \mathbb{C}^{s+1} that is nonincreasing along solutions to the recurrence relation is not new. In fact, it is easily seen to be a consequence of the Kreiss matrix Theorem, see [SW97]. However, the important point here is that we can construct a multiplier that yields directly the "energy boundedness" (or decay). The fact that the coefficients of this multiplier are integer multiples of the coefficients of P will be crucial in the analysis of Section 3, see also Proposition 2 below.

Proof. We borrow some ideas from [Gär56, Lemme 1.1] and introduce the interpolation polynomials:

$$\forall k = 1, \dots, s+1, \quad P_k(X) := a \prod_{j \neq k} (X - x_j),$$

where x_1, \dots, x_{s+1} denote the roots of P , and $a \neq 0$ its dominant coefficient. Since the roots of P are pairwise distinct, the P_k 's form a basis of $\mathbb{C}_s[X]$ and they depend in a \mathcal{C}^∞ way on the coefficients of P . We have

$$P' = \sum_{k=1}^{s+1} P_k.$$

We then consider a sequence $(v^n)_{n \in \mathbb{N}}$ with values in \mathbb{C} and compute

$$\begin{aligned} 2 \operatorname{Re} \left(\overline{\mathbf{T}(P'(\mathbf{T}) v^n)} P(\mathbf{T}) v^n \right) &= (s+1) |P(\mathbf{T}) v^n|^2 \\ &= \sum_{k=1}^{s+1} \overline{\mathbf{T}(P_k(\mathbf{T})) v^n} (\mathbf{T} - x_k) P_k(\mathbf{T}) v^n + \mathbf{T} (P_k(\mathbf{T}) v^n) (\mathbf{T} - \overline{x_k}) \overline{P_k(\mathbf{T}) v^n} \\ &\quad - \sum_{k=1}^{s+1} (\mathbf{T} - \overline{x_k}) (\overline{P_k(\mathbf{T}) v^n}) (\mathbf{T} - x_k) (P_k(\mathbf{T}) v^n) \\ &= \sum_{k=1}^{s+1} (\mathbf{T} - |x_k|^2) |P_k(\mathbf{T}) v^n|^2. \end{aligned}$$

The conclusion follows by defining:

$$\forall (w^0, \dots, w^s) \in \mathbb{C}^{s+1}, \quad q_e(w^0, \dots, w^s) := \sum_{k=1}^{s+1} |P_k(\mathbf{T}) w^0|^2, \quad (20)$$

$$q_d(w^0, \dots, w^s) := \sum_{k=1}^{s+1} (1 - |x_k|^2) |P_k(\mathbf{T}) w^0|^2. \quad (21)$$

The form q_e is positive definite because the P_k 's form a basis of $\mathbb{C}_s[X]$. The form q_d is nonnegative because the roots of P are located in $\overline{\mathbb{D}}$. Both forms depend in a \mathcal{C}^∞ way on the coefficients of P because the roots of P are simple. \square

Lemma 1 shows that the polynomial P' yields the good multiplier $\mathbf{T} P'(\mathbf{T}) v^n$ for the recurrence relation $P(\mathbf{T}) v^n = 0$. Of course, P' is not the only possible choice, though it will be our favorite one in what follows. As in [Gär56, Lemme 1.1], any polynomial of the form⁷

$$Q := \sum_{k=1}^{s+1} \alpha_k P_k, \quad \alpha_1, \dots, \alpha_{s+1} > 0,$$

provides with an energy balance of the form

$$2 \operatorname{Re} \left(\overline{\mathbf{T}(Q(\mathbf{T}) v^n)} P(\mathbf{T}) v^n \right) = (\alpha_1 + \dots + \alpha_{s+1}) |P(\mathbf{T}) v^n|^2 + (\mathbf{T} - I) (q_e(v^n, \dots, v^{n+s})) + q_d(v^n, \dots, v^{n+s}),$$

with suitable Hermitian forms q_e, q_d that have the same properties as stated in Lemma 1.

⁷The sign condition here on the coefficients α_k is the analogue of the separation condition for the roots in [Ler53, Gär56].

2.2 The energy-dissipation balance for finite difference schemes

In this Paragraph, we consider the numerical scheme (16). We introduce the following notation:

$$L := \sum_{\sigma=0}^{s+1} \mathbf{T}^\sigma Q_\sigma, \quad M := \sum_{\sigma=0}^{s+1} \sigma \mathbf{T}^\sigma Q_\sigma. \quad (22)$$

Thanks to Fourier analysis, Lemma 1 easily gives the following result:

Proposition 2 (The energy-dissipation balance law). *Let Assumptions 1 and 2 be satisfied. Then there exist a continuous coercive quadratic form E_0 and a continuous nonnegative quadratic form D_0 on $\ell^2(\mathbb{Z}^d; \mathbb{R})^{s+1}$ such that for all sequences $(v^n)_{n \in \mathbb{N}}$ with values in $\ell^2(\mathbb{Z}^d; \mathbb{R})$ and for all $n \in \mathbb{N}$, there holds*

$$2 \langle M v^n, L v^n \rangle_{-\infty, +\infty} = (s+1) \|L v^n\|_{-\infty, +\infty}^2 + (\mathbf{T} - I) E_0(v^n, \dots, v^{n+s}) + D_0(v^n, \dots, v^{n+s}).$$

In particular, for all initial data $f^0, \dots, f^s \in \ell^2(\mathbb{Z}^d; \mathbb{R})$, the solution to (16) satisfies

$$\sup_{n \in \mathbb{N}} E_0(v^n, \dots, v^{n+s}) \leq E_0(f^0, \dots, f^s),$$

and (16) is (ℓ^2) -stable.

Proof. We use the same notation v^n for the sequence $(v_j^n)_{j \in \mathbb{Z}^d}$ and the corresponding step function on \mathbb{R}^d whose value on the cell $[j_1 \Delta x_1, (j_1 + 1) \Delta x_1) \times \dots \times [j_d \Delta x_d, (j_d + 1) \Delta x_d)$ equals v_j^n . Then Plancherel Theorem gives

$$\begin{aligned} 2 \langle M v^n, L v^n \rangle_{-\infty, +\infty} - (s+1) \|L v^n\|_{-\infty, +\infty}^2 \\ = \int_{\mathbb{R}^d} 2 \operatorname{Re} \left(\overline{\mathbf{T} (P'_\zeta(\mathbf{T}) \widehat{v^n}(\xi))} P_\zeta(\mathbf{T}) \widehat{v^n}(\xi) \right) - (s+1) |P_\zeta(\mathbf{T}) \widehat{v^n}(\xi)|^2 \frac{d\xi}{(2\pi)^d}, \end{aligned}$$

where $\widehat{v^n}$ denotes the Fourier transform of v^n , and where we have let

$$P_\zeta(z) := \sum_{\sigma=0}^{s+1} \widehat{Q_\sigma} (e^{i\zeta_1}, \dots, e^{i\zeta_d}) z^\sigma, \quad \zeta_j := \xi_j \Delta x_j,$$

and $P'_\zeta(z)$ denotes the derivative of P_ζ with respect to z .

From Assumption 2, we know that for all $\zeta \in \mathbb{R}^d$, P_ζ has degree $s+1$ and has $s+1$ simple roots in $\overline{\mathbb{D}}$. We can apply Lemma 1 and get

$$\begin{aligned} 2 \langle M v^n, L v^n \rangle_{-\infty, +\infty} - (s+1) \|L v^n\|_{-\infty, +\infty}^2 \\ = \int_{\mathbb{R}^d} (\mathbf{T} - I) q_{e,\zeta}(\widehat{v^n}(\xi), \dots, \widehat{v^{n+s}}(\xi)) + q_{d,\zeta}(\widehat{v^n}(\xi), \dots, \widehat{v^{n+s}}(\xi)) \frac{d\xi}{(2\pi)^d}, \end{aligned}$$

where $q_{e,\zeta}, q_{d,\zeta}$ depend in a \mathcal{C}^∞ way on $\zeta \in \mathbb{R}^d$ and are 2π -periodic in each ζ_j . Furthermore, $q_{e,\zeta}$ is positive definite and $q_{d,\zeta}$ is nonnegative. The conclusion of Proposition 2 follows by a standard compactness argument for showing coercivity of E_0 . \square

2.3 Examples

The first basic example corresponds to the case $s = 0$ for which the multiplier provided by Proposition 2 is $Q_1 v_j^{n+1}$. In that case, the energy E_0 reads $\|Q_1 v\|_{-\infty, +\infty}^2$ (recall that Q_1 is an isomorphism) and the energy-dissipation balance law is nothing but the trivial identity

$$2 \langle Q_1 v^{n+1}, Q_1 v^{n+1} + Q_0 v^n \rangle_{-\infty, +\infty} = \|Q_1 v^{n+1} + Q_0 v^n\|_{-\infty, +\infty}^2 + \|Q_1 v^{n+1}\|_{-\infty, +\infty}^2 - \|Q_0 v^n\|_{-\infty, +\infty}^2.$$

The second line of this algebraic identity can be rewritten as

$$\|Q_1 v^{n+1}\|_{-\infty, +\infty}^2 - \|Q_1 v^n\|_{-\infty, +\infty}^2 + \|Q_1 v^n\|_{-\infty, +\infty}^2 - \|Q_0 v^n\|_{-\infty, +\infty}^2,$$

and ℓ^2 -stability for the Cauchy problem amounts to assuming that the operator norm of $Q_1^{-1} Q_0$ is not larger than 1. Hence the dissipation term $\|Q_1 v^n\|_{-\infty, +\infty}^2 - \|Q_0 v^n\|_{-\infty, +\infty}^2$ is nonnegative.

Let us now consider the leap-frog scheme in one space dimension, for which we have $s = 1$ and

$$L = \mathbf{T}^2 + \lambda a \mathbf{T} (\mathbf{S} - \mathbf{S}^{-1}) - I.$$

The corresponding dispersion relation (8) reduces to

$$z^2 + 2i\lambda a \sin \xi z - 1 = 0.$$

For $\lambda|a| < 1$, the latter equation has two simple roots $x_1(\xi), x_2(\xi)$ of modulus 1. Following the previous analysis, see (20)-(21), the form $q_{e,\zeta}$ is given by

$$q_{e,\zeta}(w^0, w^1) = |w^1 - x_1(\zeta) w^0|^2 + |w^1 - x_2(\zeta) w^0|^2 = 2|w^0|^2 + 2|w^1|^2 + 4\lambda a \operatorname{Re}(i \sin \zeta \overline{w^1} w^0),$$

and $q_{d,\zeta}$ is zero. The associated forms in Proposition 2 are $D_0 \equiv 0$ and (recall here $d = 1$):

$$E_0(v^0, v^1) = 2 \sum_{j \in \mathbb{Z}} \Delta x (v_j^0)^2 + 2 \sum_{j \in \mathbb{Z}} \Delta x (v_j^1)^2 + 2\lambda a \sum_{j \in \mathbb{Z}} \Delta x (v_{j+1}^0 - v_{j-1}^0) v_j^1.$$

The latter energy functional E_0 is coercive under the condition $\lambda|a| < 1$, which is the necessary and sufficient condition of stability for the leap-frog scheme, and E_0 is conserved for solutions to the leap-frog scheme. The conservation of E_0 is usually proved by starting from the recurrence relation

$$\forall j \in \mathbb{Z}, \quad \forall n \in \mathbb{N}, \quad u_j^{n+2} + \lambda a (u_{j+1}^{n+1} - u_{j-1}^{n+1}) - u_j^n = 0,$$

using the multiplier $u_j^{n+2} + u_j^n$, and summing with respect to j . This is equivalent, for solutions to the leap-frog scheme, to what we propose here, since our multiplier reads

$$M u_j^n = 2 u_j^{n+2} + \lambda a (u_{j+1}^{n+1} - u_{j-1}^{n+1}) = u_j^{n+2} + u_j^n + \underbrace{L u_j^n}_{=0}.$$

However, it will appear more clearly in Section 3 why our choice for $M u_j^n$ has a major advantage when considering initial boundary value problems.

Let us observe here that the energy functional E_0 is associated with a local energy density

$$E_{0,j}(v^0, v^1) := 2(v_j^0)^2 + 2(v_j^1)^2 + 2\lambda a (v_{j+1}^0 - v_{j-1}^0) v_j^1.$$

This is very specific to the leap-frog scheme. In general, the coefficients of the Hermitian forms $q_{e,\zeta}, q_{d,\zeta}$ are not trigonometric polynomials of ζ and therefore E_0, D_0 do not necessarily admit local densities. This is one main difference with [Ler53, Gär56].

3 Semigroup estimates for fully discrete initial boundary value problems

We now turn to the proof of Theorem 1 for which we shall use the results of Section 2 as a toolbox. By linearity of (5), it is sufficient to prove Theorem 1 separately in the case $(f_j^0) = \dots = (f_j^s) = 0$, and in the case $(F_j^n) = 0, (g_j^n) = 0$. The latter case is the most difficult and requires the introduction of an auxiliary set of “dissipative” boundary conditions. Solutions to (5) are always assumed to be real valued, which means that the data are real valued. For complex valued initial data and/or forcing terms, one just uses the linearity of (5).

3.1 The case with zero initial data

We first assume $(f_j^0) = \dots = (f_j^s) = 0$. By strong stability, we already know that (7) holds with a constant C that is independent of $\gamma > 0$ and $\Delta t \in (0, 1]$. Therefore, proving Theorem 1 amounts to showing the existence of a constant C , that is independent of $\gamma > 0$ and $\Delta t \in (0, 1]$ such that the solution to (5) with $(f_j^0) = \dots = (f_j^s) = 0$ satisfies

$$\sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 \leq C \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (23)$$

We thus consider a parameter $\gamma > 0$ and a time step $\Delta t \in (0, 1]$, and focus on the numerical scheme (5) with zero initial data (that is, $(f_j^0) = \dots = (f_j^s) = 0$). For all $n \in \mathbb{N}$, we extend the sequence (u_j^n) by zero for $j_1 \leq -r_1$:

$$v_j^n := \begin{cases} u_j^n & \text{if } j_1 \geq 1 - r_1, \\ 0 & \text{otherwise.} \end{cases}$$

We use Proposition 2 and compute:

$$(\mathbf{T} - I) E_0(v^n, \dots, v^{n+s}) + D_0(v^n, \dots, v^{n+s}) = 2 \langle M v^n, L v^n \rangle_{-\infty, +\infty} - (s+1) \|L v^n\|_{-\infty, +\infty}^2.$$

Due to the form of the operator L , see (22), and the fact that v_j^n vanishes for $j_1 \leq -r_1$, there holds:

$$L v_j^n = \begin{cases} \Delta t F_j^{n+s+1} & \text{if } j_1 \geq 1, \\ 0 & \text{if } j_1 \leq -r_1 - p_1, \end{cases}$$

and we thus get

$$\begin{aligned} & (\mathbf{T} - I) E_0(v^n, \dots, v^{n+s}) + D_0(v^n, \dots, v^{n+s}) \\ &= \left(\prod_{k=1}^d \Delta x_k \right) \sum_{j_1 \geq 1} \sum_{j' \in \mathbb{Z}^{d-1}} 2 \Delta t (M v_j^n) F_j^{n+s+1} - (s+1) \Delta t^2 (F_j^{n+s+1})^2 \\ &+ \left(\prod_{k=1}^d \Delta x_k \right) \sum_{j_1=1-r_1-p_1}^0 \sum_{j' \in \mathbb{Z}^{d-1}} 2 (M v_j^n) L v_j^n - (s+1) (L v_j^n)^2. \end{aligned}$$

We multiply the latter equality by $\exp(-2\gamma(n+s+1)\Delta t)$, sum with respect to n from 0 to some N and use the fact that D_0 is nonnegative. Recalling that the initial data in (5) vanish, we get

$$\underbrace{e^{-2\gamma(N+s+1)\Delta t} E_0(v^{N+1}, \dots, v^{N+s+1}) + (1 - e^{-2\gamma\Delta t}) \sum_{n=1}^N e^{-2\gamma(n+s)\Delta t} E_0(v^n, \dots, v^{n+s})}_{\geq 0} \leq S_{1,N} + S_{2,N}, \quad (24)$$

with

$$S_{1,N} := \sum_{n=0}^N e^{-2\gamma(n+s+1)\Delta t} \left(2\Delta t \langle M v^n, F^{n+s+1} \rangle_{1,+\infty} - (s+1) \Delta t^2 \|F^{n+s+1}\|_{1,+\infty}^2 \right), \quad (25)$$

and

$$S_{2,N} := \left(\prod_{k=1}^d \Delta x_k \right) \sum_{n=0}^N e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1=1-r_1-p_1}^0 \sum_{j' \in \mathbb{Z}^{d-1}} 2(M v_j^n) L v_j^n - (s+1) (L v_j^n)^2. \quad (26)$$

Let us now estimate the two source terms $S_{1,N}, S_{2,N}$ in (24). We begin with the term $S_{2,N}$ defined in (26). Let us recall that the ratio $\Delta t/\Delta x_1$ is fixed. Furthermore, the form of the operators L and M in (22) gives the estimate (recall that v_j^n vanishes for $j_1 \leq -r_1$):

$$S_{2,N} \leq C \Delta t \left(\prod_{k=2}^d \Delta x_k \right) \sum_{n=0}^N e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1=1-r_1}^{p_1} \sum_{j' \in \mathbb{Z}^{d-1}} (u_j^n)^2 + \dots + (u_j^{n+s+1})^2,$$

for a constant C that does not depend on N , γ nor on Δt . We thus have, uniformly with respect to $N \in \mathbb{N}$, $\gamma > 0$ and $\Delta t \in (0, 1]$:

$$\begin{aligned} S_{2,N} &\leq C \sum_{n=s+1}^{N+s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ &\leq C \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1,+\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}, \end{aligned}$$

where we have used the trace estimate (7) that follows from the strong stability assumption.

Let us now focus on the term $S_{1,N}$ in (24), see the defining equation (25). We use the Cauchy-Schwarz inequality and derive (using now the interior estimate in (7) that follows from the strong stability

assumption):

$$\begin{aligned}
S_{1,N} &\leq 2 \sum_{n=0}^N \Delta t e^{-2\gamma(n+s+1)\Delta t} \|M v^n\|_{1,+ \infty} \|F^{n+s+1}\|_{1,+ \infty} \\
&\leq C \sum_{n=0}^N \Delta t e^{-2\gamma(n+s+1)\Delta t} \left(\|v^{n+1}\|_{1-r_1,+ \infty} + \dots + \|v^{n+s+1}\|_{1-r_1,+ \infty} \right) \|F^{n+s+1}\|_{1,+ \infty} \\
&\leq C \frac{\gamma}{\gamma \Delta t + 1} \sum_{n=s+1}^{N+s+1} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1,+ \infty}^2 + C \frac{\gamma \Delta t + 1}{\gamma} \sum_{n=s+1}^{N+s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1,+ \infty}^2 \\
&\leq C \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1,+ \infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1,\cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}.
\end{aligned}$$

Ignoring the nonnegative term on the left hand-side of (24) and using the coercivity of E_0 , we have proved that there exists a constant $C > 0$ that is uniform with respect to $N, \gamma, \Delta t$ such that:

$$\begin{aligned}
e^{-2\gamma(N+s+1)\Delta t} \|v^{N+s+1}\|_{-\infty,+ \infty}^2 &\leq C \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1,+ \infty}^2 \right. \\
&\quad \left. + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|g_{j_1,\cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\},
\end{aligned}$$

which yields (23) and therefore the validity of Theorem 1 in the case of zero initial data.

3.2 Construction of dissipative boundary conditions

In this paragraph, we consider an auxiliary problem for which we shall be able to prove simultaneously an optimal semigroup estimate and a trace estimate for the solution. More precisely, we shall prove the following result.

Theorem 2. *Let Assumptions 1, 2 and 3 be satisfied. Then for all $P_1 \in \mathbb{N}$, there exists a constant $C_{P_1} > 0$ such that, for all initial data $(f_j^0), \dots, (f_j^s) \in \ell^2(\mathbb{Z}^d)$ and for all source term $(g_j^n)_{j_1 \leq 0, n \geq s+1}$ that satisfies*

$$\forall \Gamma > 0, \quad \sum_{n \geq s+1} e^{-2\Gamma n} \sum_{j_1 \leq 0} \|g_{j_1,\cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 < +\infty,$$

there exists a unique sequence $(u_j^n)_{j \in \mathbb{Z}^d, n \in \mathbb{N}}$ solution to

$$\begin{cases} L u_j^n = 0, & j \in \mathbb{Z}^d, \quad j_1 \geq 1, \quad n \geq 0, \\ M u_j^n = g_j^{n+s+1}, & j \in \mathbb{Z}^d, \quad j_1 \leq 0, \quad n \geq 0, \\ u_j^n = f_j^n, & j \in \mathbb{Z}^d, \quad n = 0, \dots, s. \end{cases} \quad (27)$$

Moreover for all $\gamma > 0$ and $\Delta t \in (0, 1]$, this solution satisfies

$$\begin{aligned} \sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{-\infty, +\infty}^2 &+ \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{-\infty, +\infty}^2 + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{P_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ &\leq C_{P_1} \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (28) \end{aligned}$$

Theorem 2 justifies why we advocate the choice $M u_j^n = 2 u_j^{n+2} + \lambda a(u_{j+1}^{n+1} - u_{j-1}^{n+1})$ rather than the more standard $u_j^{n+2} + u_j^n$ as a multiplier for the leap-frog scheme. Despite repeated efforts, we have not been able to prove the estimate (28) when using the numerical boundary condition $u_j^{n+2} + u_j^n$ on $j_1 \leq 0$, in conjunction with the leap-frog scheme on $j_1 \geq 1$.

Proof. Let us first quickly observe that the solution to (27) is well-defined since, as long as we have determined the solution up to a time index $n + s$, $n \geq 0$, then u^{n+s+1} is sought as a solution to an equation of the form

$$Q_{s+1} u^{n+s+1} = F,$$

where F belongs to $\ell^2(\mathbb{Z}^d)$ (this is due to the form of L and M , see (22)). Hence u^n is uniquely defined and belongs to $\ell^2(\mathbb{Z}^d)$ for all $n \in \mathbb{N}$.

The proof of Theorem 2 starts again with the application of Proposition 2. Using the nonnegativity of the dissipation form D_0 , we get⁸

$$(\mathbf{T} - I) E_0(u^n, \dots, u^{n+s}) + (s+1) \|L u^n\|_{-\infty, +\infty}^2 \leq 2 \langle M u^n, L u^n \rangle_{-\infty, +\infty} = 2 \langle g^{n+s+1}, L u^n \rangle_{-\infty, 0}.$$

By the Young inequality, we get

$$(\mathbf{T} - I) E_0(u^n, \dots, u^{n+s}) + \frac{s+1}{2} \|L u^n\|_{-\infty, +\infty}^2 \leq \frac{2}{s+1} \|g^{n+s+1}\|_{-\infty, 0}^2.$$

We multiply the latter inequality by $\exp(-2\gamma(n+s+1)\Delta t)$, sum from $n = 0$ to some arbitrary N and already derive the estimate (here we use again the fact that $\Delta t/\Delta x_1$ is a fixed positive constant):

$$\begin{aligned} \sup_{n \geq 1} e^{-2\gamma(n+s)\Delta t} E_0(u^n, \dots, u^{n+s}) &+ (1 - e^{-2\gamma\Delta t}) \sum_{n \geq 0} e^{-2\gamma(n+s)\Delta t} E_0(u^n, \dots, u^{n+s}) \\ &+ \sum_{n \geq 0} \Delta t e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1 \in \mathbb{Z}} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ &\leq C \left\{ e^{-2\gamma s \Delta t} E_0(f^0, \dots, f^s) + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \end{aligned}$$

Using the coercivity of E_0 and the inequality

$$1 - e^{-2\gamma\Delta t} \geq \frac{\gamma\Delta t}{\gamma\Delta t + 1},$$

⁸Since $L u_j^n = 0$ for $j_1 \geq 1$, one could also write $\|L u^n\|_{-\infty, 0}^2$ rather than $\|L u^n\|_{-\infty, +\infty}^2$ on the left hand-side of the inequality.

we have therefore derived the estimate

$$\begin{aligned} & \sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{-\infty, +\infty}^2 + \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{-\infty, +\infty}^2 \\ & + \sum_{n \geq 0} \Delta t e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1 \in \mathbb{Z}} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ & \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}, \quad (29) \end{aligned}$$

where the constant C is independent of γ , Δt and on the solution (u_j^n) . In order to prove (28), the main remaining task is to derive the trace estimate for (u_j^n) . This is done by first dealing with the case where $\gamma \Delta t$ is large.

• From the definition of the operator L , see (22), there exists a constant $C > 0$ and an integer J such that

$$(L u_j^n)^2 \geq \frac{1}{2} (Q_{s+1} u_j^{n+s+1})^2 - C \sum_{\sigma=0}^s \sum_{|\ell| \leq J} (u_{j+\ell}^{n+\sigma})^2.$$

Since Q_{s+1} is an isomorphism, there exists a constant $c > 0$ such that

$$\sum_{j \in \mathbb{Z}^d} (L u_j^n)^2 \geq c \sum_{j \in \mathbb{Z}^d} (u_j^{n+s+1})^2 - \frac{1}{c} \sum_{\sigma=0}^s \sum_{j \in \mathbb{Z}^d} (u_j^{n+\sigma})^2.$$

Multiplying by $\exp(-2\gamma(n+s+1)\Delta t)$ and summing with respect to $n \in \mathbb{N}$, we get

$$\begin{aligned} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \in \mathbb{Z}} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 & \leq C \left\{ \sum_{n \geq 0} \Delta t e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1 \in \mathbb{Z}} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right. \\ & \quad \left. + e^{-2\gamma \Delta t} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \in \mathbb{Z}} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \end{aligned}$$

Choosing $\gamma \Delta t$ large enough, that is $\gamma \Delta t \geq \ln R_0$ for some numerical constant $R_0 > 1$ that depends only on the (fixed) coefficients of the operator L , we have derived the estimate

$$\begin{aligned} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \in \mathbb{Z}} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 & \leq C \left\{ \sum_{n \geq 0} \Delta t e^{-2\gamma(n+s+1)\Delta t} \sum_{j_1 \in \mathbb{Z}} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right. \\ & \quad \left. + e^{-2\gamma \Delta t} \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 \right\}. \end{aligned}$$

It remains to use (29) and we get an even better estimate than (28) which we were originally aiming at:

$$\sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \in \mathbb{Z}} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}.$$

This gives a control of infinitely many traces and not only finitely many (this restriction to finitely many traces will appear in the regime where $\gamma \Delta t$ can be small).

• From now on, we have fixed a constant $R_0 > 1$ such that (28) holds for $\gamma \Delta t \geq \ln R_0$ and we thus assume $\gamma \Delta t \in (0, \ln R_0]$. We also know that the estimate (29) holds, independently of the value of $\gamma \Delta t$, and we now wish to estimate the traces of the solution (u_j^n) for finitely many values of j_1 .

We first observe from (29) that for all $\gamma > 0$ and $\Delta t \in (0, 1]$, there exists a constant $C_{\gamma, \Delta t}$ such that

$$\forall n \in \mathbb{N}, \quad e^{-2\gamma n \Delta t} \sum_{j \in \mathbb{Z}^d} (u_j^n)^2 \leq C_{\gamma, \Delta t}.$$

In particular, for any $j_1 \in \mathbb{Z}$, the Laplace-Fourier transforms $\widehat{u_{j_1}}$ of the step functions

$$u_{j_1} : (t, y) \in \mathbb{R}^+ \times \mathbb{R}^{d-1} \mapsto u_j^n \quad \text{if } (t, y) \in [n \Delta t, (n+1) \Delta t) \times \prod_{k=2}^d [j_k \Delta x_k, (j_k+1) \Delta x_k),$$

is well-defined on $\{\tau \in \mathbb{C}, \operatorname{Re} \tau > 0\} \times \mathbb{R}^{d-1}$. The dual variables are denoted $\tau = \gamma + i\theta$, $\gamma > 0$, and $\eta = (\eta_2, \dots, \eta_d) \in \mathbb{R}^{d-1}$. It will also be convenient to introduce the notation $\eta_\Delta := (\eta_2 \Delta x_2, \dots, \eta_d \Delta x_d)$. Given $\Gamma > 0$, the sequence $(\widehat{u_{j_1}}(\Gamma + i\theta, \eta))_{j_1 \in \mathbb{Z}}$ belongs to $\ell^2(\mathbb{Z})$ for almost every $(\theta, \eta) \in \mathbb{R} \times \mathbb{R}^{d-1}$.

We first show the following estimate.

Lemma 2. *With $R_0 > 1$ fixed as above, there exists a constant $C > 0$ such that for all $\gamma > 0$ and $\Delta t \in (0, 1]$ satisfying $\gamma \Delta t \in (0, \ln R_0]$, there holds*

$$\begin{aligned} & \sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} \left| \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1} (e^{(\gamma+i\theta)\Delta t}, \eta_\Delta) \widehat{u_{j_1+\ell_1}}(\gamma + i\theta, \eta) \right|^2 d\theta d\eta \\ & + \sum_{j_1 \leq 0} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} \left| \sum_{\ell_1=-r_1}^{p_1} e^{(\gamma+i\theta)\Delta t} \partial_z a_{\ell_1} (e^{(\gamma+i\theta)\Delta t}, \eta_\Delta) \widehat{u_{j_1+\ell_1}}(\gamma + i\theta, \eta) \right|^2 d\theta d\eta \\ & \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (30) \end{aligned}$$

Proof of Lemma 2. Given $\tau = \gamma + i\theta$ and η , we compute (here $j_1 \in \mathbb{Z}$ is fixed):

$$\sum_{\ell_1=-r_1}^{p_1} a_{\ell_1} (e^{\tau \Delta t}, \eta_\Delta) \widehat{u_{j_1+\ell_1}}(\tau, \eta) = \widehat{L u_{j_1, \cdot}}(\tau, \eta) + \frac{1 - e^{-\tau \Delta t}}{\tau} \sum_{\sigma=1}^{s+1} \sum_{\sigma'=0}^{\sigma-1} e^{(\sigma-\sigma')\tau \Delta t} \mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta), \quad (31)$$

$$\sum_{\ell_1=-r_1}^{p_1} e^{\tau \Delta t} \partial_z a_{\ell_1} (e^{\tau \Delta t}, \eta_\Delta) \widehat{u_{j_1+\ell_1}}(\tau, \eta) = \widehat{M u_{j_1, \cdot}}(\tau, \eta) + \frac{1 - e^{-\tau \Delta t}}{\tau} \sum_{\sigma=1}^{s+1} \sum_{\sigma'=0}^{\sigma-1} \sigma e^{(\sigma-\sigma')\tau \Delta t} \mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta). \quad (32)$$

where, in (31) and (32), we have set

$$\mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta) = \sum_{\ell_1=-r_1}^{p_1} \left(\sum_{\ell'=-r'}^{p'} a_{\ell, \sigma} e^{i\ell' \cdot \eta_\Delta} \right) \widehat{f_{j_1+\ell_1, \cdot}^{\sigma'}}(\eta),$$

which corresponds to the partial Fourier transform with respect to $y = (x_2, \dots, x_d) \in \mathbb{R}^{d-1}$, of the step function associated with the sequence $(Q_\sigma f_j^{\sigma'})$ (no Laplace transform here).

We need to estimate integrals with respect to (θ, η) of the right hand side of (31) and (32). The first term on the right of (31) and (32) are easy. For instance, we have (applying Plancherel Theorem):

$$\begin{aligned} \sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\widehat{L u_{j_1, \cdot}}(\tau, \eta)|^2 d\theta d\eta &= (2\pi)^d \sum_{j_1 \in \mathbb{Z}} \sum_{n \geq 0} \int_{n\Delta t}^{(n+1)\Delta t} e^{-2\gamma s} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 ds \\ &= (2\pi)^d \frac{1 - e^{-2\gamma \Delta t}}{2\gamma \Delta t} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \in \mathbb{Z}} \|L u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2. \end{aligned}$$

We now recall that $\gamma \Delta t$ is restricted to the interval $(0, \ln R_0]$, and we use (29) to derive

$$\sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\widehat{L u_{j_1, \cdot}}(\tau, \eta)|^2 d\theta d\eta \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}.$$

Similarly, we have

$$\sum_{j_1 \leq 0} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\widehat{M u_{j_1, \cdot}}(\tau, \eta)|^2 d\theta d\eta = (2\pi)^d \frac{1 - e^{-2\gamma \Delta t}}{2\gamma \Delta t} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|M u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2,$$

which we can again uniformly estimate by the right hand side of (30).

Going back to the right hand side terms in (31) and (32), we find that there only remains for proving (30) to estimate the integral (here there are finitely many values of σ and σ'):

$$\sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} \left| \frac{1 - e^{-\tau \Delta t}}{\tau} \right|^2 |\mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta)|^2 d\theta d\eta = \left(\int_{\mathbb{R}} \left| \frac{1 - e^{-\tau \Delta t}}{\tau} \right|^2 d\theta \right) \sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R}^{d-1}} |\mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta)|^2 d\eta,$$

where we have applied Fubini Theorem. Applying first Plancherel Theorem with respect to the $d-1$ last space variables, we get

$$\sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R}^{d-1}} |\mathcal{F}_{j_1}^{\sigma, \sigma'}(\eta)|^2 d\eta \leq C \sum_{j_1 \in \mathbb{Z}} \sum_{j' \in \mathbb{Z}^{d-1}} \left(\prod_{k=2}^d \Delta x_k \right) (f_j^{\sigma'})^2 \leq \frac{C}{\Delta t} \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2.$$

The conclusion then follows by computing

$$\int_{\mathbb{R}} \left| \frac{1 - e^{-\tau \Delta t}}{\tau} \right|^2 d\theta = 2\pi \Delta t \frac{1 - e^{-2\gamma \Delta t}}{2\gamma \Delta t},$$

and by recalling that $\gamma \Delta t$ belongs to $(0, \ln R_0]$. We can eventually bound the integrals on the left hand side of (30) by estimating separately the integrals of each term on the right hand side of (31) and (32). \square

The conclusion now relies on the following crucial result.

Lemma 3 (The trace estimate). *Let Assumptions 1, 2 and 3 be satisfied. Let $R_0 > 1$ be fixed as above and let $P_1 \in \mathbb{N}$. Then there exists a constant $C_{P_1} > 0$ such that for all $z \in \mathcal{U}$ with $|z| \leq R_0$, for all $\eta \in \mathbb{R}^{d-1}$ and for all sequence $(w_{j_1})_{j_1 \in \mathbb{Z}} \in \ell^2(\mathbb{Z}; \mathbb{C})$, there holds*

$$\sum_{j_1 = -r_1 - p_1}^{P_1} |w_{j_1}|^2 \leq C_{P_1} \left\{ \sum_{j_1 \in \mathbb{Z}} \left| \sum_{\ell_1 = -r_1}^{p_1} a_{\ell_1}(z, \eta) w_{j_1 + \ell_1} \right|^2 + \sum_{j_1 \leq 0} \left| \sum_{\ell_1 = -r_1}^{p_1} z \partial_z a_{\ell_1}(z, \eta) w_{j_1 + \ell_1} \right|^2 \right\}. \quad (33)$$

Recall that the functions a_{ℓ_1} , $\ell_1 = -r_1, \dots, p_1$, are defined in (9).

The proof of Lemma 3 is rather long. Before giving it in full details, we indicate how Lemma 3 yields the result of Theorem 2. We apply Lemma 3 to $z = \exp(\tau \Delta t)$, $\tau = \gamma + i\theta$ with $\gamma \Delta t \in (0, \ln R_0]$, $\eta_\Delta \in \mathbb{R}^{d-1}$ and the sequence $(\widehat{u}_{j_1}(\tau, \eta))_{j_1 \in \mathbb{Z}}$. We then integrate (33) with respect to (θ, η) and use Lemma 2 to derive

$$\sum_{j_1=-r_1-p_1}^{P_1} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\widehat{u}_{j_1}(\gamma + i\theta, \eta)|^2 d\theta d\eta \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}.$$

It remains to apply Plancherel Theorem and we get

$$\begin{aligned} \sum_{j_1=-r_1-p_1}^{P_1} \sum_{n \in \mathbb{N}} \frac{1 - e^{-2\gamma \Delta t}}{2\gamma \Delta t} \Delta t e^{-2\gamma n \Delta t} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \end{aligned}$$

Recalling that $\gamma \Delta t$ is restricted to the interval $(0, \ln R_0]$, we have thus derived the trace estimate

$$\sum_{n \in \mathbb{N}} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=-r_1-p_1}^{P_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1 \leq 0} \|g_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}.$$

Combined with the semigroup and interior estimate (29), this gives the estimate (28) of Theorem 2 for $\gamma \Delta t \in (0, \ln R_0]$. \square

Proof of Lemma 3. Let us recall that the functions a_{ℓ_1} are 2π -periodic with respect to each coordinate of η . We can therefore restrict to $\eta \in [0, 2\pi]^{d-1}$ rather than considering $\eta \in \mathbb{R}^{d-1}$. We argue by contradiction and assume that the conclusion to Lemma 3 does not hold. This means the following, up to normalizing and extracting subsequences; there exist three sequences (indexed by $k \in \mathbb{N}$):

- a sequence $(w^k)_{k \in \mathbb{N}}$ with values in $\ell^2(\mathbb{Z}; \mathbb{C})$ such that $(w_{-r_1-p_1}^k, \dots, w_{P_1}^k)$ belongs to the unit sphere of $\mathbb{C}^{P_1+r_1+p_1+1}$ for all k , and $(w_{-r_1-p_1}^k, \dots, w_{P_1}^k)$ converges towards $(\underline{w}_{-r_1-p_1}, \dots, \underline{w}_{P_1})$ as k tends to infinity,
- a sequence $(z^k)_{k \in \mathbb{N}}$ with values in $\mathcal{U} \cap \{\zeta \in \mathbb{C}, |\zeta| \leq R_0\}$, which converges towards $\underline{z} \in \overline{\mathcal{U}}$,
- a sequence $(\eta^k)_{k \in \mathbb{N}}$ with values in $[0, 2\pi]^{d-1}$, which converges towards $\underline{\eta} \in [0, 2\pi]^{d-1}$,

and these sequences satisfy:

$$\lim_{k \rightarrow +\infty} \sum_{j_1 \in \mathbb{Z}} \left| \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z^k, \eta^k) w_{j_1+\ell_1}^k \right|^2 + \sum_{j_1 \leq 0} \left| \sum_{\ell_1=-r_1}^{p_1} z^k \partial_z a_{\ell_1}(z^k, \eta^k) w_{j_1+\ell_1}^k \right|^2 = 0. \quad (34)$$

We are going to show that (34) implies that $(\underline{w}_{-r_1-p_1}, \dots, \underline{w}_{P_1})$ must be zero, which will yield a contradiction since this vector must have norm 1.

• Let us first show that each component $(w_{j_1}^k)_{k \in \mathbb{N}}$, $j_1 \in \mathbb{Z}$, has a limit as k tends to infinity. This is already clear for $j_1 = -r_1 - p_1, \dots, P_1$. For $j_1 > P_1$, we argue by induction. From (34), we have

$$\lim_{k \rightarrow +\infty} \sum_{\ell_1 = -r_1}^{p_1} a_{\ell_1}(z^k, \eta^k) w_{P_1 - p_1 + 1 + \ell_1}^k = 0,$$

and by Assumption 3, we know that $a_{p_1}(\underline{z}, \underline{\eta})$ is nonzero. Hence $(w_{P_1+1}^k)_{k \in \mathbb{N}}$ converges towards

$$-\frac{1}{a_{p_1}(\underline{z}, \underline{\eta})} \sum_{\ell_1 = -r_1}^{p_1-1} a_{\ell_1}(\underline{z}, \underline{\eta}) \underline{w}_{P_1 - p_1 + 1 + \ell_1},$$

which we define as \underline{w}_{P_1+1} . We can argue by induction in the same way for all indices $j_1 > P_1 + 1$, but also for indices $j_1 < -r_1 - p_1$ because the function a_{-r_1} also does not vanish on $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$.

Using (34), we have thus shown that for each $j_1 \in \mathbb{Z}$, $(w_{j_1}^k)_{k \in \mathbb{N}}$ tends towards some limit \underline{w}_{j_1} as k tends to infinity, and the sequence \underline{w} , which does not necessarily belong to $\ell^2(\mathbb{Z}; \mathbb{C})$, satisfies the induction relations:

$$\forall j_1 \in \mathbb{Z}, \quad \sum_{\ell_1 = -r_1}^{p_1} a_{\ell_1}(\underline{z}, \underline{\eta}) \underline{w}_{j_1 + \ell_1} = 0, \quad (35)$$

$$\forall j_1 \leq 0, \quad \sum_{\ell_1 = -r_1}^{p_1} \underline{z} \partial_z a_{\ell_1}(\underline{z}, \underline{\eta}) \underline{w}_{j_1 + \ell_1} = 0. \quad (36)$$

• The induction relation (35) is the one that arises in [GKS72, Mic83] and all the works that deal with strong stability. The main novelty here is to use simultaneously (35) for controlling the unstable components of $(\underline{w}_{-r_1-p_1}, \dots, \underline{w}_{-1})$ and (36) for controlling the stable components of $(\underline{w}_{-r_1-p_1}, \dots, \underline{w}_{-1})$. The fact that \underline{w} satisfies simultaneously (35) and (36) for $j_1 \leq 0$ automatically annihilates the central components. This sketch of proof is made precise below.

We define the source terms:

$$F_{j_1}^k := \sum_{\ell_1 = -r_1}^{p_1} a_{\ell_1}(z^k, \eta^k) w_{j_1 + \ell_1}^k, \quad G_{j_1}^k := \sum_{\ell_1 = -r_1}^{p_1} z^k \partial_z a_{\ell_1}(z^k, \eta^k) w_{j_1 + \ell_1}^k,$$

which, according to (34), satisfy

$$\lim_{k \rightarrow 0} \sum_{j_1 \in \mathbb{Z}} |F_{j_1}^k|^2 = 0, \quad \lim_{k \rightarrow 0} \sum_{j_1 \leq 0} |G_{j_1}^k|^2 = 0. \quad (37)$$

We also introduce the vectors (here T denotes transposition)

$$W_{j_1}^k := \left(w_{j_1+p_1}^k, \dots, w_{j_1+1-r_1}^k \right)^T, \quad \underline{W}_{j_1} := \left(\underline{w}_{j_1+p_1}, \dots, \underline{w}_{j_1+1-r_1} \right)^T,$$

and the matrices:

$$\mathbb{L}(z, \eta) := \begin{pmatrix} -a_{p_1-1}(z, \eta)/a_{p_1}(z, \eta) & \cdots & \cdots & -a_{-r_1}(z, \eta)/a_{p_1}(z, \eta) \\ 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ 0 & 0 & 1 & 0 \end{pmatrix} \in \mathcal{M}_{p_1+r_1}(\mathbb{C}), \quad (38)$$

$$\mathbb{M}(z, \eta) := \begin{pmatrix} -\partial_z a_{p_1-1}(z, \eta)/\partial_z a_{p_1}(z, \eta) & \cdots & \cdots & -\partial_z a_{-r_1}(z, \eta)/\partial_z a_{p_1}(z, \eta) \\ 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ 0 & 0 & 1 & 0 \end{pmatrix} \in \mathcal{M}_{p_1+r_1}(\mathbb{C}). \quad (39)$$

The matrix \mathbb{L} is well-defined on $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$ according to Assumption 3. The matrix \mathbb{M} is also well-defined on $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$ because for any $\eta \in \mathbb{R}^{d-1}$, Assumption 3 asserts that $a_{p_1}(\cdot, \eta)$ is a nonconstant polynomial whose roots lie in \mathbb{D} . From the Gauss-Lucas Theorem, the roots of $\partial_z a_{p_1}(\cdot, \eta)$ lie in the convex hull of those of $a_{p_1}(\cdot, \eta)$. Therefore $\partial_z a_{p_1}(\cdot, \eta)$ does not vanish on $\overline{\mathcal{U}}$. In the same way, $\partial_z a_{-r_1}(\cdot, \eta)$ does not vanish on $\overline{\mathcal{U}}$.

With our above notation, the vectors $W_{j_1}^k, \underline{W}_{j_1}$, satisfy the one step induction relations:

$$\forall j_1 \in \mathbb{Z}, \quad W_{j_1+1}^k = \mathbb{L}(z^k, \eta^k) W_{j_1}^k + \left(F_{j_1+1}^k/a_{p_1}(z^k, \eta^k), 0, \dots, 0 \right)^T, \quad \underline{W}_{j_1+1} = \mathbb{L}(z, \underline{\eta}) \underline{W}_{j_1}, \quad (40)$$

$$\forall j_1 \leq -1, \quad W_{j_1+1}^k = \mathbb{M}(z^k, \eta^k) W_{j_1}^k + \left(G_{j_1+1}^k/(z^k \partial_z a_{p_1}(z^k, \eta^k)), 0, \dots, 0 \right)^T, \quad \underline{W}_{j_1+1} = \mathbb{M}(z, \underline{\eta}) \underline{W}_{j_1}. \quad (41)$$

• From Assumption 3 and the above application of the Gauss-Lucas Theorem, we already know that both matrices $\mathbb{L}(z, \eta)$ and $\mathbb{M}(z, \eta)$ are invertible for $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$. Furthermore, Assumption 2 shows that $\mathbb{L}(z, \eta)$ has no eigenvalue on \mathbb{S}^1 for $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$. This property dates back at least to [Kre68]. However, central eigenvalues on \mathbb{S}^1 may occur for \mathbb{L} when z belongs to \mathbb{S}^1 . The crucial point for proving Lemma 3 is that Assumption 2 precludes central eigenvalues of \mathbb{M} for all $z \in \overline{\mathcal{U}}$. Namely, for all $z \in \overline{\mathcal{U}}$ and all $\eta \in \mathbb{R}^{d-1}$, $\mathbb{M}(z, \eta)$ has no eigenvalue on \mathbb{S}^1 . This property holds because otherwise, for some $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$, there would exist a solution $\kappa_1 \in \mathbb{S}^1$ to the dispersion relation

$$\sum_{\ell_1=-r_1}^{p_1} z \partial_z a_{\ell_1}(z, \eta) \kappa_1^{\ell_1} = 0.$$

For convenience, the coordinates of η are denoted (η_2, \dots, η_d) . Using the definition (9) of a_{ℓ_1} , and defining $\kappa := (\kappa_1, e^{i\eta_2}, \dots, e^{i\eta_d})$, we have found a root $z \in \overline{\mathcal{U}}$ to the relation

$$\sum_{\sigma=1}^{s+1} \sigma \widehat{Q}_\sigma(\kappa) z^{\sigma-1} = 0, \quad (42)$$

but this is not possible because the $s+1$ roots (in z) to the dispersion relation (8) are simple and belong to $\overline{\mathbb{D}}$. The Gauss-Lucas Theorem thus shows that the roots to the relation (42) belong to \mathbb{D} (and therefore not to $\overline{\mathcal{U}}$).

At this stage, we know that the eigenvalues of $\mathbb{M}(z, \eta)$, $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$, split into two groups: those in \mathcal{U} , which we call the unstable ones, and those in \mathbb{D} , which we call the stable ones. For $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$,

we then introduce the spectral projector $\Pi_{\mathbb{M}}^s(z, \eta)$, resp. $\Pi_{\mathbb{M}}^u(z, \eta)$, of $\mathbb{M}(z, \eta)$ on the generalized eigenspace associated with eigenvalues in \mathbb{D} , resp. \mathcal{U} . We can then integrate the first induction relation in (41) and get

$$\Pi_{\mathbb{M}}^s(z^k, \eta^k) W_0^k = \frac{1}{z^k \partial_z a_{p_1}(z^k, \eta^k)} \sum_{j_1 \leq 0} \mathbb{M}(z^k, \eta^k)^{|j_1|} \Pi_{\mathbb{M}}^s(z^k, \eta^k) \left(G_{j_1}^k, 0, \dots, 0 \right)^T.$$

The projector $\Pi_{\mathbb{M}}^s$ depends continuously on $(z, \eta) \in \overline{\mathcal{U}} \times \mathbb{R}^{d-1}$. Furthermore, since the spectrum of \mathbb{M} does not meet \mathbb{S}^1 even for $z \in \mathbb{S}^1$, there exists a constant $C > 0$ and a $\delta \in (0, 1)$ that are independent of $k \in \mathbb{N}$ and such that

$$\forall j_1 \leq 0, \quad \|\mathbb{M}(z^k, \eta^k)^{|j_1|} \Pi_{\mathbb{M}}^s(z^k, \eta^k)\| \leq C \delta^{|j_1|}.$$

We thus get a uniform estimate with respect to k :

$$|\Pi_{\mathbb{M}}^s(z^k, \eta^k) W_0^k|^2 \leq C \sum_{j_1 \leq 0} |G_{j_1}^k|^2.$$

Passing to the limit and using (37), we get $\Pi_{\mathbb{M}}^s(\underline{z}, \underline{\eta}) \underline{W}_0 = 0$, or in other words $\underline{W}_0 = \Pi_{\mathbb{M}}^u(\underline{z}, \underline{\eta}) \underline{W}_0$.

- The sequence $(\underline{W}_{j_1})_{j_1 \leq 0}$ satisfies both induction relations (40) and (41). Due to the form of the companion matrices \mathbb{L} and \mathbb{M} , see (38)-(39), we can conclude that the vector \underline{W}_0 belongs to the generalized eigenspace (of either \mathbb{L} or \mathbb{M}) associated with the common eigenvalues of $\mathbb{M}(\underline{z}, \underline{\eta})$ and $\mathbb{L}(\underline{z}, \underline{\eta})$. We have already seen that $\mathbb{M}(\underline{z}, \underline{\eta})$ has no eigenvalue on \mathbb{S}^1 and $\underline{W}_0 = \Pi_{\mathbb{M}}^u(\underline{z}, \underline{\eta}) \underline{W}_0$, so we can conclude that \underline{W}_0 belongs to the generalized eigenspace of \mathbb{L} associated with those common eigenvalues of $\mathbb{M}(\underline{z}, \underline{\eta})$ and $\mathbb{L}(\underline{z}, \underline{\eta})$ in \mathcal{U} .

The matrix $\mathbb{L}(\underline{z}, \underline{\eta})$ has N^u eigenvalues in \mathcal{U} , N^s in \mathbb{D} and N^c on \mathbb{S}^1 . (Since \underline{z} may belong to \mathbb{S}^1 , N^c is not necessarily zero.) With obvious notations, we let $\Pi_{\mathbb{L}}^{u,s,c}(z, \eta)$ denote the corresponding spectral projectors of \mathbb{L} for (z, η) sufficiently close to $(\underline{z}, \underline{\eta})$. In particular, the eigenvalues corresponding to $\Pi_{\mathbb{L}}^u(z, \eta)$ lie in \mathcal{U} uniformly away from \mathbb{S}^1 for (z, η) sufficiently close to $(\underline{z}, \underline{\eta})$. We can then integrate the first induction relation in (40) and derive (for k sufficiently large):

$$\Pi_{\mathbb{L}}^u(z^k, \eta^k) W_0^k = -\frac{1}{a_{p_1}(z^k, \eta^k)} \sum_{j_1 \geq 0} \mathbb{L}(z^k, \eta^k)^{-j_1-1} \Pi_{\mathbb{L}}^u(z^k, \eta^k) \left(F_{j_1}^k, 0, \dots, 0 \right)^T.$$

Using the uniform exponential decay of $\mathbb{L}(z^k, \eta^k)^{-j_1-1} \Pi_{\mathbb{L}}^u(z^k, \eta^k)$ and (37), we finally end up with

$$\Pi_{\mathbb{L}}^u(\underline{z}, \underline{\eta}) \underline{W}_0 = 0.$$

Since \underline{W}_0 belongs to the generalized eigenspace of \mathbb{L} associated with those common eigenvalues of $\mathbb{M}(\underline{z}, \underline{\eta})$ and $\mathbb{L}(\underline{z}, \underline{\eta})$ in \mathcal{U} , we can conclude that \underline{W}_0 equals zero. Applying the induction relation (40), the whole sequence $(\underline{W}_{j_1})_{j_1 \in \mathbb{Z}}$ is zero which yields the expected contradiction. \square

The crucial property that we use in the proof of Lemma 3 is the fact that up to $z \in \mathbb{S}^1$, the eigenvalues of $\mathbb{M}(z, \eta)$ lie either in \mathbb{D} or \mathcal{U} . For the leap-frog scheme, this property would not be true if we had imposed the auxiliary numerical boundary condition $u_j^{n+2} + u_j^n$ rather than $2u_j^{n+2} + \lambda a(u_{j+1}^{n+1} - u_{j-1}^{n+1})$.

Let us also observe that we have used the fact that a_{p_1} and a_{-r_1} are nonconstant in order to study the induction relation (36). There might be some schemes for which a_{p_1} and/or a_{-r_1} are constant but for which one can still apply similar arguments as in the previous proof, even though (36) is an induction relation with fewer steps than (35). In this respect, Assumption 3 might be relaxed in specific applications.

Remark 2. The auxiliary problem (27) is in general not of the same form as (5) because in (27) one has to impose infinitely many numerical boundary conditions. This is due to the fact that the stencil of M incorporates points “on the left” with respect to the first space variable. A remarkable exception occurs for explicit schemes with $s = 0$, for in that case the multiplier $M v_j^n$ reads v_j^{n+1} and (27) is exactly the auxiliary problem considered (and labeled (2.7)) in [CG11] where one imposes Dirichlet boundary conditions on finitely many boundary meshes (just use $g_j^n = 0$ for $j_1 \leq -r_1$). The reader can then check that the energy-dissipation balance law of Proposition 2 in that case ($s = 0$, $Q_1 = I$) coincides exactly with the algebra involved in the derivation of the estimate (2.12) in [CG11]. The reader can also check that for $s = 0$, and $Q_1 = I$, Lemma 3 becomes a rather trivial exercise...

There still remains the problem of constructing a set of dissipative numerical boundary conditions of the same form as (5) with $s \geq 1$, that is with finitely many numerical boundary conditions, and for which one can prove by hand both a semigroup and a trace estimate as in Theorem 2.

3.3 End of the proof

As explained in the introduction of Section 3, the linearity of (5) reduces the proof of Theorem 1 to the case $(F_j^n) = 0$, $(g_j^n) = 0$, since we have already dealt with the case of zero initial data. We thus focus on (5) with $(F_j^n) = 0$ and $(g_j^n) = 0$, and write the corresponding solution (u_j^n) as $u_j^n = v_j^n + w_j^n$, where the sequence (v_j^n) solves:

$$\begin{cases} L v_j^n = 0, & j \in \mathbb{Z}^d, \quad j_1 \geq 1, \quad n \geq 0, \\ M v_j^n = 0, & j \in \mathbb{Z}^d, \quad j_1 \leq 0, \quad n \geq 0, \\ v_j^n = f_j^n, & j \in \mathbb{Z}^d, \quad n = 0, \dots, s, \end{cases} \quad (43)$$

and (w_j^n) solves:

$$\begin{cases} L w_j^n = 0, & j \in \mathbb{Z}^d, \quad j_1 \geq 1, \quad n \geq 0, \\ w_j^{n+s+1} + \sum_{\sigma=0}^{s+1} B_{j_1, \sigma} w_{1, j'}^{n+\sigma} = \tilde{g}_j^{n+s+1}, & j \in \mathbb{Z}^d, \quad j_1 = 1 - r_1, \dots, 0, \quad n \geq 0, \\ w_j^n = 0, & j \in \mathbb{Z}^d, \quad n = 0, \dots, s. \end{cases} \quad (44)$$

For $v_j^n + w_j^n$ to coincide with the solution (u_j^n) to (5), it is sufficient to extend the initial data f_j^0, \dots, f_j^s by zero for $j_1 \leq -r_1$, which provides with the initial data in (43) on all \mathbb{Z}^d , and to define the boundary source term in (44) by:

$$\tilde{g}_j^{n+s+1} := -v_j^{n+s+1} - \sum_{\sigma=0}^{s+1} B_{j_1, \sigma} v_{1, j'}^{n+\sigma}. \quad (45)$$

We can estimate the solution (v_j^n) to (43) by applying Theorem 2. In particular, the trace estimate:

$$\sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{P_1} \|v_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2,$$

for $P_1 = \max(p_1, q_1 + 1)$ gives (recall the definition (45) of \tilde{g}_j^{n+s+1}):

$$\begin{aligned} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|\tilde{g}_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 &\leq C \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{\max(p_1, q_1+1)} \|v_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ &\leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2. \end{aligned}$$

We can apply Theorem 1 to the solution (w_j^n) to (44) because the initial data in (44) vanish. We get:

$$\begin{aligned} \sup_{n \geq 0} e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2 + \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2 \\ + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|w_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^0 \|\tilde{g}_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ \leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2. \end{aligned}$$

Combining with the similar estimate provided by Theorem 2 for (v_j^n) , we end up with the expected estimate:

$$\begin{aligned} \sup_{n \geq 0} e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 + \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 \\ + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|u_{j_1, \cdot}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2, \end{aligned}$$

which completes the proof of Theorem 1.

4 Conclusion and perspectives

Let us first observe that in [Wad90], WADE has constructed symmetrizers for deriving stability estimates for multistep schemes, even in the case of variable coefficients. His conditions for constructing a symmetrizer are less restrictive than Assumption 2. However, the symmetrizer in [Wad90] is genuinely nonlocal and it is therefore not clear that it may be useful for boundary value problems. The main novelty here is to construct a *local* multiplier whose properties allow for the design of an auxiliary *dissipative* boundary value problem. This is our key to Theorem 1.

In this article we have always discarded the dissipation term provided by the nonnegative form D_0 . For the approximation of parabolic equations, this term may give some extra dissipation, but a crucial point to keep in mind is that the coefficients of the numerical scheme are assumed to be constant (which may in turn yield rather severe CFL conditions for implicit approximations of parabolic equations). Hence it does not seem very clear that our approach will yield stability estimates with “optimal” CFL conditions when approximating parabolic equations. This extension is left to further study in the future.

The main possible improvement of Theorem 1 would consist of assuming that only the roots to (8) that lie on \mathbb{S}^1 are simple. Here we have assumed that all the roots, including those in \mathbb{D} are simple. If

we could manage to deal with multiple roots in \mathbb{D} , then Theorem 1 would be applicable to numerical approximations of the transport equation (12) that are based on Adams-Bashforth methods of order 3 or higher (such methods have 0 as a root of multiplicity 2 or more at the zero frequency).

The results in this paper achieve the proof of a "weak form" of the conjecture in [KW93] that strong stability, in the sense of Definition 1, implies semigroup stability. However, an even stronger assumption was made in [KW93], namely that the sole fulfillment of the interior estimate

$$\frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 \leq C \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2,$$

when both the initial *and boundary data* for (5) vanish, does imply semigroup stability. The analogous conjecture for partial differential equations seems to be still open so far, but we do hope that our multiplier technique may yield some insight for dealing with the strong form of the conjecture in [KW93].

Acknowledgments This article was completed while the author was visiting the Institute of Mathematical Science at Nanjing University. The author warmly thanks Professor Yin Huicheng and the Institute for their hospitality during this visit.

References

- [AG76] S. Abarbanel and D. Gottlieb. A note on the leap-frog scheme in two and three space dimensions. *J. Computational Phys.*, 21(3):351–355, 1976.
- [AG79] S. Abarbanel and D. Gottlieb. Stability of two-dimensional initial boundary value problems using leap-frog type schemes. *Math. Comp.*, 33(148):1145–1155, 1979.
- [BGS07] S. Benzoni-Gavage and D. Serre. *Multidimensional hyperbolic partial differential equations*. Oxford University Press, 2007. First-order systems and applications.
- [CG11] J.-F. Coulombel and A. Gloria. Semigroup stability of finite difference schemes for multidimensional hyperbolic initial boundary value problems. *Math. Comp.*, 80(273):165–203, 2011.
- [Cou09] J.-F. Coulombel. Stability of finite difference schemes for hyperbolic initial boundary value problems. *SIAM J. Numer. Anal.*, 47(4):2844–2871, 2009.
- [Cou13] J.-F. Coulombel. Stability of finite difference schemes for hyperbolic initial boundary value problems. In *HCDTE Lecture Notes. Part I. Nonlinear Hyperbolic PDEs, Dispersive and Transport Equations*, pages 97–225. American Institute of Mathematical Sciences, 2013.
- [Cou14] J.-F. Coulombel. Fully discrete hyperbolic initial boundary value problems with nonzero initial data. *Preprint*, <http://arxiv.org/abs/1412.0851>, 2014.
- [Går56] L. Gårding. Solution directe du problème de Cauchy pour les équations hyperboliques. In *La théorie des équations aux dérivées partielles*, Colloques Internationaux du C. N. R. S., pages 71–90. C. N. R. S., Paris, 1956.
- [GKO95] B. Gustafsson, H.-O. Kreiss, and J. Oliger. *Time dependent problems and difference methods*. John Wiley & Sons, 1995.

- [GKS72] B. Gustafsson, H.-O. Kreiss, and A. Sundström. Stability theory of difference approximations for mixed initial boundary value problems. II. *Math. Comp.*, 26(119):649–686, 1972.
- [GT81] M. Goldberg and E. Tadmor. Scheme-independent stability criteria for difference approximations of hyperbolic initial-boundary value problems. II. *Math. Comp.*, 36(154):603–626, 1981.
- [HNW93] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving ordinary differential equations. I*. Springer-Verlag, second edition, 1993. Nonstiff problems.
- [HW96] E. Hairer and G. Wanner. *Solving ordinary differential equations. II*. Springer-Verlag, second edition, 1996. Stiff and differential-algebraic problems.
- [Kre68] H.-O. Kreiss. Stability theory for difference approximations of mixed initial boundary value problems. I. *Math. Comp.*, 22:703–714, 1968.
- [KW93] H.-O. Kreiss and L. Wu. On the stability definition of difference approximations for the initial-boundary value problem. *Appl. Numer. Math.*, 12(1-3):213–227, 1993.
- [Ler53] J. Leray. *Hyperbolic differential equations*. The Institute for Advanced Study, Princeton, N. J., 1953.
- [Mic83] D. Michelson. Stability theory of difference approximations for multidimensional initial-boundary value problems. *Math. Comp.*, 40(161):1–45, 1983.
- [Oli74] J. Oliger. Fourth order difference methods for the initial boundary-value problem for hyperbolic equations. *Math. Comp.*, 28:15–25, 1974.
- [Rau72] J. Rauch. \mathcal{L}^2 is a continuable initial condition for Kreiss’ mixed problems. *Comm. Pure Appl. Math.*, 25:265–285, 1972.
- [RM67] R. D. Richtmyer and K. W. Morton. *Difference methods for initial value problems*. Graduate Texts in Mathematics. Interscience Publishers John Wiley & Sons, 1967. Theory and applications.
- [Slo83] D. M. Sloan. Boundary conditions for a fourth order hyperbolic difference scheme. *Math. Comp.*, 41:1–11, 1983.
- [SW97] J. C. Strikwerda and B. A. Wade. A survey of the Kreiss matrix theorem for power bounded families of matrices and its extensions. In *Linear operators (Warsaw, 1994)*, volume 38 of *Banach Center Publ.*, pages 339–360. Polish Acad. Sci., 1997.
- [TE05] L. N. Trefethen and M. Embree. *Spectra and pseudospectra*. Princeton University Press, 2005. The behavior of nonnormal matrices and operators.
- [Tho72] J. M. Thomas. Discrétisation des conditions aux limites dans les schémas saute-mouton. *Rev. Française Automat. Informat. Recherche Opérationnelle*, 6(Ser. R-2):31–44, 1972.
- [Tre84] L. N. Trefethen. Instability of difference models for hyperbolic initial boundary value problems. *Comm. Pure Appl. Math.*, 37:329–367, 1984.
- [Wad90] B. A. Wade. Symmetrizable finite difference operators. *Math. Comp.*, 54(190):525–543, 1990.

- [Wu95] L. Wu. The semigroup stability of the difference approximations for initial-boundary value problems. *Math. Comp.*, 64(209):71–88, 1995.